

# King County Renovation: Most Effective Ways to Increase the Value of Your Home

**Author:** E. Berke Tezcan

---

## TABLE OF CONTENTS

*Click to jump to matching Markdown Header.*

- [INTRODUCTION](#)
  - [OBTAIN](#)
  - [SCRUB/EXPLORE](#)
  - [MODEL](#)
  - [INTERPRET](#)
  - [CONCLUSIONS/RECOMMENDATIONS](#)
- 

## INTRODUCTION

For this project, we were hired by a home owner in King County, Washington who wants to renovate their home. They would like us to analyze the real estate data of the county and give them insights as to what to focus their renovation efforts on in order to increase their property's value.

Real estate prices are affected by a myriad of -what we can define as- internal parameters like the square footage, floor count, bedroom count, finishes, how many cars can fit into the garage etc. There are also parameters that we can define as external and can not (easily) be changed. These include attributes like zipcode, latitude, longitude, view from the house, and school districts. In order to accurately model and pinpoint the most important parameters that affect the sale price of a home, we need to incorporate both internal and external parameters.

We are given a dataset that includes information about the real estate in King County and will be using this dataset to create a multiple linear regression (MLR) model. We defined our goal to be to find the top 3 internal parameters that affect a home's sale price the most in King County specifically. This will ensure that the home owner can actually keep these parameters in mind while renovating rather than getting insights about external parameters that they can't necessarily do anything to change.

## OBTAIN

### Data Understanding/EDA

In [1]:

```
import pandas as pd
import seaborn as sns
```

```

import matplotlib.pyplot as plt
# from sklearn.linear_model import LinearRegression
from sklearn.preprocessing import OneHotEncoder
import statsmodels.api as sm
import statsmodels.stats.api as sms
import statsmodels.formula.api as smf
import numpy as np
from scipy import stats

```

In [2]: df = pd.read\_csv('data/kc\_house\_data.csv')

In [3]: pd.set\_option('display.max\_columns',0)  
df.head()

Out[3]:

|   | id         | date       | price    | bedrooms | bathrooms | sqft_living | sqft_lot | floors | waterfront | vi  |
|---|------------|------------|----------|----------|-----------|-------------|----------|--------|------------|-----|
| 0 | 7129300520 | 10/13/2014 | 221900.0 |          | 3         | 1.00        | 1180     | 5650   | 1.0        | NaN |
| 1 | 6414100192 | 12/9/2014  | 538000.0 |          | 3         | 2.25        | 2570     | 7242   | 2.0        | 0.0 |
| 2 | 5631500400 | 2/25/2015  | 180000.0 |          | 2         | 1.00        | 770      | 10000  | 1.0        | 0.0 |
| 3 | 2487200875 | 12/9/2014  | 604000.0 |          | 4         | 3.00        | 1960     | 5000   | 1.0        | 0.0 |
| 4 | 1954400510 | 2/18/2015  | 510000.0 |          | 3         | 2.00        | 1680     | 8080   | 1.0        | 0.0 |

In [4]: df.info()

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 21597 entries, 0 to 21596
Data columns (total 21 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   id          21597 non-null   int64  
 1   date         21597 non-null   object  
 2   price        21597 non-null   float64 
 3   bedrooms     21597 non-null   int64  
 4   bathrooms    21597 non-null   float64 
 5   sqft_living  21597 non-null   int64  
 6   sqft_lot     21597 non-null   int64  
 7   floors       21597 non-null   float64 
 8   waterfront   19221 non-null   float64 
 9   view         21534 non-null   float64 
 10  condition    21597 non-null   int64  
 11  grade        21597 non-null   int64  
 12  sqft_above   21597 non-null   int64  
 13  sqft_basement 21597 non-null   object  
 14  yr_built     21597 non-null   int64  
 15  yr_renovated 17755 non-null   float64 
 16  zipcode      21597 non-null   int64  
 17  lat          21597 non-null   float64 
 18  long         21597 non-null   float64 
 19  sqft_living15 21597 non-null   int64  
 20  sqft_lot15   21597 non-null   int64  
dtypes: float64(8), int64(11), object(2)
memory usage: 3.5+ MB

```

In [5]: #Looking at unique values within the column to see if there are any out-of-place looking  
df['bathrooms'].unique()

Out[5]: array([1. , 2.25, 3. , 2. , 4.5 , 1.5 , 2.5 , 1.75, 2.75, 3.25, 4. ,

```
3.5 , 0.75, 4.75, 5. , 4.25, 3.75, 1.25, 5.25, 6. , 0.5 , 5.5 ,
6.75, 5.75, 8. , 7.5 , 7.75, 6.25, 6.5 ])
```

In [6]: `#Looking at unique values within the column to see if there are any out-of-place lookin  
df['condition'].unique()`

Out[6]: `array([3, 5, 4, 1, 2], dtype=int64)`

In [7]: `#Looking at unique values within the column to see if there are any out-of-place lookin  
df['grade'].unique()`

Out[7]: `array([ 7, 6, 8, 11, 9, 5, 10, 12, 4, 3, 13], dtype=int64)`

In [8]: `#Looking at unique values within the column to see if there are any out-of-place lookin  
df['view'].unique()`

Out[8]: `array([ 0., nan, 3., 4., 2., 1.])`

In [9]: `#Checking to see how many unique homes we have in the dataset.  
pd.DataFrame(df['id'].unique())`

Out[9]:

|              | <b>0</b>   |
|--------------|------------|
| <b>0</b>     | 7129300520 |
| <b>1</b>     | 6414100192 |
| <b>2</b>     | 5631500400 |
| <b>3</b>     | 2487200875 |
| <b>4</b>     | 1954400510 |
| ...          | ...        |
| <b>21415</b> | 263000018  |
| <b>21416</b> | 6600060120 |
| <b>21417</b> | 1523300141 |
| <b>21418</b> | 291310100  |
| <b>21419</b> | 1523300157 |

21420 rows × 1 columns

We have 21,597 individual home sales associated with 21,420 unique homes.

## SCRUB/EXPLORE

### Checking for and Addressing Null Values

In [10]: `df.isna().sum()`

|              |   |
|--------------|---|
| <b>id</b>    | 0 |
| <b>date</b>  | 0 |
| <b>price</b> | 0 |

```

bedrooms          0
bathrooms         0
sqft_living       0
sqft_lot          0
floors            0
waterfront        2376
view              63
condition         0
grade              0
sqft_above         0
sqft_basement      0
yr_built           0
yr_renovated      3842
zipcode            0
lat                0
long               0
sqft_living15     0
sqft_lot15         0
dtype: int64

```

We have 3 columns with null values in them (waterfront, view and yr\_renovated as shown above) which we will need to address.

## view Column - Null Values

```
In [11]: df['view'].value_counts(dropna=False)
```

```

Out[11]: 0.0    19422
2.0     957
3.0     508
1.0     330
4.0     317
NaN      63
Name: view, dtype: int64

```

We can see that the majority (approx. 90%) of the data points have zeros for their view column value. Therefore, to address the missing values in the view column we can take the more conservative approach and say that these 63 houses/apartments did not have a view (had a value of 0) with minimal impact to the overall dataset since we have 21597 data points overall.

```
In [12]: df['view'].fillna(0, inplace=True)
df.isna().sum()
```

```

Out[12]: id          0
date         0
price         0
bedrooms      0
bathrooms     0
sqft_living    0
sqft_lot       0
floors         0
waterfront    2376
view          0
condition      0
grade          0
sqft_above      0
sqft_basement   0
yr_built        0
yr_renovated    3842
zipcode         0
lat             0
long            0

```

```
sqft_living15      0
sqft_lot15        0
dtype: int64
```

## yr\_renovated Column - Null Values

Prior to checking for null values, we should check to see if there are inconsistencies that we need to address between the yr\_built and yr\_renovated columns. For example, if both columns are filled in with the same year, that may indicate that one of the columns is incorrect and may need to be replaced with a 0. Another example would be if the yr\_built column showed a larger year compared to the yr\_renovated since this would suggest that the building was built after it was renovated which doesn't make sense.

```
In [13]: df[df['yr_built']==df['yr_renovated']]
```

```
Out[13]:   id  date  price  bedrooms  bathrooms  sqft_living  sqft_lot  floors  waterfront  view  condition  grac
```

```
In [14]: df[(df['yr_built']>df['yr_renovated']) & (df['yr_renovated']!=0)]
```

```
Out[14]:   id  date  price  bedrooms  bathrooms  sqft_living  sqft_lot  floors  waterfront  view  condition  grac
```

We found no inconsistencies between the yr\_built and yr\_renovated columns and can carry on with addressing the null values in the yr\_renovated column.

```
In [15]: df['yr_renovated'].unique()
```

```
Out[15]: array([  0., 1991.,    nan, 2002., 2010., 1992., 2013., 1994., 1978.,
       2005., 2003., 1984., 1954., 2014., 2011., 1983., 1945., 1990.,
       1988., 1977., 1981., 1995., 2000., 1999., 1998., 1970., 1989.,
       2004., 1986., 2007., 1987., 2006., 1985., 2001., 1980., 1971.,
       1979., 1997., 1950., 1969., 1948., 2009., 2015., 1974., 2008.,
       1968., 2012., 1963., 1951., 1962., 1953., 1993., 1996., 1955.,
       1982., 1956., 1940., 1976., 1946., 1975., 1964., 1973., 1957.,
       1959., 1960., 1967., 1965., 1934., 1972., 1944., 1958.])
```

```
In [16]: df['yr_renovated'].value_counts(dropna=False)
```

```
Out[16]: 0.0      17011
NaN      3842
2014.0     73
2003.0     31
2013.0     31
...
1944.0      1
1948.0      1
1976.0      1
1934.0      1
1953.0      1
Name: yr_renovated, Length: 71, dtype: int64
```

Similar to the view column, we are seeing that majority (approx. 79%) of the data points have 0 for their yr\_renovated column value which we are making the assumption that the place was not renovated in this case. Chances are that the places with null values in the yr\_renovated column were also not renovated, so we can go ahead and fill the null values with zeros in this case as well.

```
In [17]: df['yr_renovated'].fillna(0, inplace=True)
```

```
In [18]: df.isna().sum()
```

```
Out[18]: id          0
date         0
price        0
bedrooms     0
bathrooms    0
sqft_living  0
sqft_lot     0
floors       0
waterfront   2376
view         0
condition    0
grade         0
sqft_above   0
sqft_basement 0
yr_built     0
yr_renovated 0
zipcode      0
lat          0
long         0
sqft_living15 0
sqft_lot15   0
dtype: int64
```

## waterfront Column - Null Values

```
In [19]: df['waterfront'].unique()
```

```
Out[19]: array([nan,  0.,  1.])
```

```
In [20]: df['waterfront'].value_counts(dropna=False)
```

```
Out[20]: 0.0    19075
NaN     2376
1.0     146
Name: waterfront, dtype: int64
```

Similar to the other two columns we discussed, majority of the homes (approx. 88%) have 0 as their waterfront column value. Therefore we can assume that the 2,376 null values can also be filled with 0.

```
In [21]: df['waterfront'].fillna(0, inplace=True)
```

```
In [22]: df.isna().sum()
```

```
Out[22]: id          0
date         0
price        0
bedrooms     0
bathrooms    0
sqft_living  0
sqft_lot     0
floors       0
waterfront   0
view         0
condition    0
grade         0
```

```

sqft_above      0
sqft_basement   0
yr_built        0
yr_renovated    0
zipcode         0
lat             0
long            0
sqft_living15   0
sqft_lot15      0
dtype: int64

```

## Addressing Placeholder Values

For some reason, unlike the other columns that have sqft information, the sqft\_basement column's data type seems to be object instead of integer. If we want to be able to accurately use this column's information we need to change the dtype to int64 similar to other sqft columns.

In [23]: `df.info()`

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 21597 entries, 0 to 21596
Data columns (total 21 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   id          21597 non-null   int64  
 1   date         21597 non-null   object  
 2   price        21597 non-null   float64 
 3   bedrooms     21597 non-null   int64  
 4   bathrooms    21597 non-null   float64 
 5   sqft_living  21597 non-null   int64  
 6   sqft_lot     21597 non-null   int64  
 7   floors       21597 non-null   float64 
 8   waterfront   21597 non-null   float64 
 9   view         21597 non-null   float64 
 10  condition    21597 non-null   int64  
 11  grade        21597 non-null   int64  
 12  sqft_above   21597 non-null   int64  
 13  sqft_basement 21597 non-null   object  
 14  yr_built     21597 non-null   int64  
 15  yr_renovated 21597 non-null   float64 
 16  zipcode      21597 non-null   int64  
 17  lat          21597 non-null   float64 
 18  long         21597 non-null   float64 
 19  sqft_living15 21597 non-null   int64  
 20  sqft_lot15   21597 non-null   int64  
dtypes: float64(8), int64(11), object(2)
memory usage: 3.5+ MB

```

In [24]: `df['sqft_basement'].dtype`

Out[24]: `dtype('O')`

In order to find the reason we can take a quick look at the unique values that this column has. We expect to see an object (or better known as string) value that is different than all the other values.

In [25]: `df['sqft_basement'].unique()`

```

Out[25]: array(['0.0', '400.0', '910.0', '1530.0', '?', '730.0', '1700.0', '300.0',
               '970.0', '760.0', '720.0', '700.0', '820.0', '780.0', '790.0',
               '330.0', '1620.0', '360.0', '588.0', '1510.0', '410.0', '990.0',
               '600.0', '560.0', '550.0', '1000.0', '1600.0', '500.0', '1040.0'],
              dtype='|S11')

```

```
'880.0', '1010.0', '240.0', '265.0', '290.0', '800.0', '540.0',
'710.0', '840.0', '380.0', '770.0', '480.0', '570.0', '1490.0',
'620.0', '1250.0', '1270.0', '120.0', '650.0', '180.0', '1130.0',
'450.0', '1640.0', '1460.0', '1020.0', '1030.0', '750.0', '640.0',
'1070.0', '490.0', '1310.0', '630.0', '2000.0', '390.0', '430.0',
'850.0', '210.0', '1430.0', '1950.0', '440.0', '220.0', '1160.0',
'860.0', '580.0', '2060.0', '1820.0', '1180.0', '200.0', '1150.0',
'1200.0', '680.0', '530.0', '1450.0', '1170.0', '1080.0', '960.0',
'280.0', '870.0', '1100.0', '460.0', '1400.0', '660.0', '1220.0',
'900.0', '420.0', '1580.0', '1380.0', '475.0', '690.0', '270.0',
'350.0', '935.0', '1370.0', '980.0', '1470.0', '160.0', '950.0',
'50.0', '740.0', '1780.0', '1900.0', '340.0', '470.0', '370.0',
'140.0', '1760.0', '130.0', '520.0', '890.0', '1110.0', '150.0',
'1720.0', '810.0', '190.0', '1290.0', '670.0', '1800.0', '1120.0',
'1810.0', '60.0', '1050.0', '940.0', '310.0', '930.0', '1390.0',
'610.0', '1830.0', '1300.0', '510.0', '1330.0', '1590.0', '920.0',
'1320.0', '1420.0', '1240.0', '1960.0', '1560.0', '2020.0',
'1190.0', '2110.0', '1280.0', '250.0', '2390.0', '1230.0', '170.0',
'830.0', '1260.0', '1410.0', '1340.0', '590.0', '1500.0', '1140.0',
'260.0', '100.0', '320.0', '1480.0', '1060.0', '1284.0', '1670.0',
'1350.0', '2570.0', '1090.0', '110.0', '2500.0', '90.0', '1940.0',
'1550.0', '2350.0', '2490.0', '1481.0', '1360.0', '1135.0',
'1520.0', '1850.0', '1660.0', '2130.0', '2600.0', '1690.0',
'243.0', '1210.0', '1024.0', '1798.0', '1610.0', '1440.0',
'1570.0', '1650.0', '704.0', '1910.0', '1630.0', '2360.0',
'1852.0', '2090.0', '2400.0', '1790.0', '2150.0', '230.0', '70.0',
'1680.0', '2100.0', '3000.0', '1870.0', '1710.0', '2030.0',
'875.0', '1540.0', '2850.0', '2170.0', '506.0', '906.0', '145.0',
'2040.0', '784.0', '1750.0', '374.0', '518.0', '2720.0', '2730.0',
'1840.0', '3480.0', '2160.0', '1920.0', '2330.0', '1860.0',
'2050.0', '4820.0', '1913.0', '80.0', '2010.0', '3260.0', '2200.0',
'415.0', '1730.0', '652.0', '2196.0', '1930.0', '515.0', '40.0',
'2080.0', '2580.0', '1548.0', '1740.0', '235.0', '861.0', '1890.0',
'2220.0', '792.0', '2070.0', '4130.0', '2250.0', '2240.0',
'1990.0', '768.0', '2550.0', '435.0', '1008.0', '2300.0', '2610.0',
'666.0', '3500.0', '172.0', '1816.0', '2190.0', '1245.0', '1525.0',
'1880.0', '862.0', '946.0', '1281.0', '414.0', '2180.0', '276.0',
'1248.0', '602.0', '516.0', '176.0', '225.0', '1275.0', '266.0',
'283.0', '65.0', '2310.0', '10.0', '1770.0', '2120.0', '295.0',
'207.0', '915.0', '556.0', '417.0', '143.0', '508.0', '2810.0',
'20.0', '274.0', '248.0'], dtype=object)
```

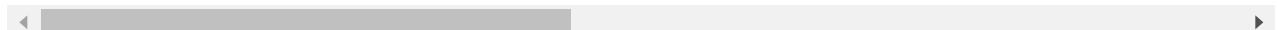
It seems like the '?' values are being used as placeholder values and are causing the sqft\_basement column to be a object type column. Let's take a quick look at how many '?' values there are.

In [26]: `df[df['sqft_basement']=='?']`

|              | <b>id</b>  | <b>date</b> | <b>price</b> | <b>bedrooms</b> | <b>bathrooms</b> | <b>sqft_living</b> | <b>sqft_lot</b> | <b>floors</b> | <b>waterfro</b> |
|--------------|------------|-------------|--------------|-----------------|------------------|--------------------|-----------------|---------------|-----------------|
| <b>6</b>     | 1321400060 | 6/27/2014   | 257500.0     | 3               | 2.25             | 1715               | 6819            | 2.0           | (               |
| <b>18</b>    | 16000397   | 12/5/2014   | 189000.0     | 2               | 1.00             | 1200               | 9850            | 1.0           | (               |
| <b>42</b>    | 7203220400 | 7/7/2014    | 861990.0     | 5               | 2.75             | 3595               | 5639            | 2.0           | (               |
| <b>79</b>    | 1531000030 | 3/23/2015   | 720000.0     | 4               | 2.50             | 3450               | 39683           | 2.0           | (               |
| <b>112</b>   | 2525310310 | 9/16/2014   | 272500.0     | 3               | 1.75             | 1540               | 12600           | 1.0           | (               |
| ...          | ...        | ...         | ...          | ...             | ...              | ...                | ...             | ...           | ...             |
| <b>21442</b> | 3226049565 | 7/11/2014   | 504600.0     | 5               | 3.00             | 2360               | 5000            | 1.0           | (               |
| <b>21447</b> | 1760650900 | 7/21/2014   | 337500.0     | 4               | 2.50             | 2330               | 4907            | 2.0           | (               |

|              | <b>id</b>  | <b>date</b> | <b>price</b> | <b>bedrooms</b> | <b>bathrooms</b> | <b>sqft_living</b> | <b>sqft_lot</b> | <b>floors</b> | <b>waterfront</b> |
|--------------|------------|-------------|--------------|-----------------|------------------|--------------------|-----------------|---------------|-------------------|
| <b>21473</b> | 6021503707 | 1/20/2015   | 352500.0     |                 | 2                | 2.50               | 980             | 1010          | 3.0               |
| <b>21519</b> | 2909310100 | 10/15/2014  | 332000.0     |                 | 4                | 2.50               | 2380            | 5737          | 2.0               |
| <b>21581</b> | 191100405  | 4/21/2015   | 1580000.0    |                 | 4                | 3.25               | 3410            | 10125         | 2.0               |

454 rows × 21 columns



In [27]: `df['sqft_basement'].value_counts()`

```
Out[27]: 0.0      12826
?
454
600.0     217
500.0     209
700.0     208
...
935.0      1
666.0      1
2570.0     1
1852.0     1
243.0      1
Name: sqft_basement, Length: 304, dtype: int64
```

It seems like only 454 rows out of 21,597 have '?' as their value. Once again we see that majority (approx. 59%) of the values in this column have a value of 0. Therefore, taking the conservative approach, we can once again assume that these placeholder values are to have 0 values and change the dtype to int.

In [28]: `df['sqft_basement'] = df['sqft_basement'].map(lambda x: x.replace('?', '0') if x=='?' else float(x))`  
`# df['sqft_basement'] = df['sqft_basement'].map(Lambda x: int(float(x)))`  
`df['sqft_basement'] = df['sqft_basement'].astype('float').astype('int64')`

In [29]: `df['sqft_basement'].dtype`

Out[29]: `dtype('int64')`

## Feature Engineering - renovated

Since the objective of our data analysis is to give homeowners 3 recommendations on how to renovate their homes, the yr\_renovated column in this dataset is not very helpful and will not result in an actionable recommendation since homeowners can't go back in time to renovate their homes during the optimal years. Instead looking at whether a house was renovated or not may reveal more and reinforce the idea that renovated homes would have an increase in their value. Therefore, we should create a feature called 'renovated' based on the 'yr\_renovated' column.

In [30]: `df['renovated'] = df['yr_renovated']!=0`  
`df['renovated'] = df['renovated'].astype('int')`  
`df.head()`

|          | <b>id</b>  | <b>date</b> | <b>price</b> | <b>bedrooms</b> | <b>bathrooms</b> | <b>sqft_living</b> | <b>sqft_lot</b> | <b>floors</b> | <b>waterfront</b> | <b>view</b> |
|----------|------------|-------------|--------------|-----------------|------------------|--------------------|-----------------|---------------|-------------------|-------------|
| <b>0</b> | 7129300520 | 10/13/2014  | 221900.0     |                 | 3                | 1.00               | 1180            | 5650          | 1.0               | 0.0         |

|          | <b>id</b>  | <b>date</b> | <b>price</b> | <b>bedrooms</b> | <b>bathrooms</b> | <b>sqft_living</b> | <b>sqft_lot</b> | <b>floors</b> | <b>waterfront</b> | <b>vi</b> |
|----------|------------|-------------|--------------|-----------------|------------------|--------------------|-----------------|---------------|-------------------|-----------|
| <b>1</b> | 6414100192 | 12/9/2014   | 538000.0     |                 | 3                | 2.25               | 2570            | 7242          | 2.0               | 0.0       |
| <b>2</b> | 5631500400 | 2/25/2015   | 180000.0     |                 | 2                | 1.00               | 770             | 10000         | 1.0               | 0.0       |
| <b>3</b> | 2487200875 | 12/9/2014   | 604000.0     |                 | 4                | 3.00               | 1960            | 5000          | 1.0               | 0.0       |
| <b>4</b> | 1954400510 | 2/18/2015   | 510000.0     |                 | 3                | 2.00               | 1680            | 8080          | 1.0               | 0.0       |

## Feature Engineering - has\_basement

We think that having a finished basement in a home may be a more important metric compared to how large the basement is in terms of sqft. Therefore we wanted to engineer a feature that would show whether the homes had a finished basement or not. During our EDA we noticed that  $\text{sqft\_living} = \text{sqft\_above} + \text{sqft\_basement}$ . This would suggest that if a basement was unfinished and was not occupiable it was not included in the total livable sqft number which is denoted by  $\text{sqft\_living}$ . Therefore, basements that were unfinished would have the same 0 value as homes without any basements. So we are assuming that the homes with 0's as their  $\text{sqft\_basement}$  value, may still have a basement but that their basement would not be an occupiable, livable space.

```
In [31]: df['has_basement'] = df['sqft_basement']!=0
df['has_basement'] = df['has_basement'].astype('int')
df.head()
```

|          | <b>id</b>  | <b>date</b> | <b>price</b> | <b>bedrooms</b> | <b>bathrooms</b> | <b>sqft_living</b> | <b>sqft_lot</b> | <b>floors</b> | <b>waterfront</b> | <b>vi</b> |
|----------|------------|-------------|--------------|-----------------|------------------|--------------------|-----------------|---------------|-------------------|-----------|
| <b>0</b> | 7129300520 | 10/13/2014  | 221900.0     |                 | 3                | 1.00               | 1180            | 5650          | 1.0               | 0.0       |
| <b>1</b> | 6414100192 | 12/9/2014   | 538000.0     |                 | 3                | 2.25               | 2570            | 7242          | 2.0               | 0.0       |
| <b>2</b> | 5631500400 | 2/25/2015   | 180000.0     |                 | 2                | 1.00               | 770             | 10000         | 1.0               | 0.0       |
| <b>3</b> | 2487200875 | 12/9/2014   | 604000.0     |                 | 4                | 3.00               | 1960            | 5000          | 1.0               | 0.0       |
| <b>4</b> | 1954400510 | 2/18/2015   | 510000.0     |                 | 3                | 2.00               | 1680            | 8080          | 1.0               | 0.0       |

## Checking for Duplicates

The best way to check for duplicated data for this dataset specifically is by using the 'id' column as each home was given a unique ID number.

```
In [32]: df[df['id'].duplicated()]
```

|            | <b>id</b>  | <b>date</b> | <b>price</b> | <b>bedrooms</b> | <b>bathrooms</b> | <b>sqft_living</b> | <b>sqft_lot</b> | <b>floors</b> | <b>waterfro</b> |   |
|------------|------------|-------------|--------------|-----------------|------------------|--------------------|-----------------|---------------|-----------------|---|
| <b>94</b>  | 6021501535 | 12/23/2014  | 700000.0     |                 | 3                | 1.50               | 1580            | 5000          | 1.0             | ( |
| <b>314</b> | 4139480200 | 12/9/2014   | 1400000.0    |                 | 4                | 3.25               | 4290            | 12103         | 1.0             | ( |
| <b>325</b> | 7520000520 | 3/11/2015   | 240500.0     |                 | 2                | 1.00               | 1240            | 12092         | 1.0             | ( |
| <b>346</b> | 3969300030 | 12/29/2014  | 239900.0     |                 | 4                | 1.00               | 1000            | 7134          | 1.0             | ( |

|       | <b>id</b>  | <b>date</b> | <b>price</b> | <b>bedrooms</b> | <b>bathrooms</b> | <b>sqft_living</b> | <b>sqft_lot</b> | <b>floors</b> | <b>waterfront</b> |
|-------|------------|-------------|--------------|-----------------|------------------|--------------------|-----------------|---------------|-------------------|
| 372   | 2231500030 | 3/24/2015   | 530000.0     | 4               | 2.25             | 2180               | 10754           | 1.0           | (                 |
| ...   | ...        | ...         | ...          | ...             | ...              | ...                | ...             | ...           | ...               |
| 20165 | 7853400250 | 2/19/2015   | 645000.0     | 4               | 3.50             | 2910               | 5260            | 2.0           | (                 |
| 20597 | 2724049222 | 12/1/2014   | 220000.0     | 2               | 2.50             | 1000               | 1092            | 2.0           | (                 |
| 20654 | 8564860270 | 3/30/2015   | 502000.0     | 4               | 2.50             | 2680               | 5539            | 2.0           | (                 |
| 20764 | 6300000226 | 5/4/2015    | 380000.0     | 4               | 1.00             | 1200               | 2171            | 1.5           | (                 |
| 21565 | 7853420110 | 5/4/2015    | 625000.0     | 3               | 3.00             | 2780               | 6000            | 2.0           | (                 |

177 rows × 23 columns

There are duplicated entries for the same houses but since these are additional sales of the same house, they are valid data points and there is no reason to drop them. So we will be leaving these duplicates in our dataset as they will give us a more complete picture of the market.

## Linearity Check

Since we were specifically asked to use an MLR model for this project, we need to address the assumptions of this model to make sure that the model is accurate. One of the assumptions for the model is that the parameters we will be using to infer price need to have a linear relationship against price.

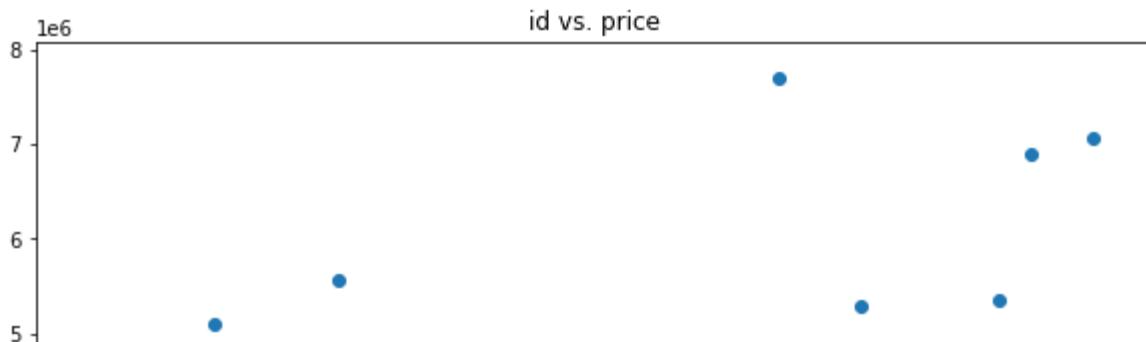
### Checking for Linearity of Parameters

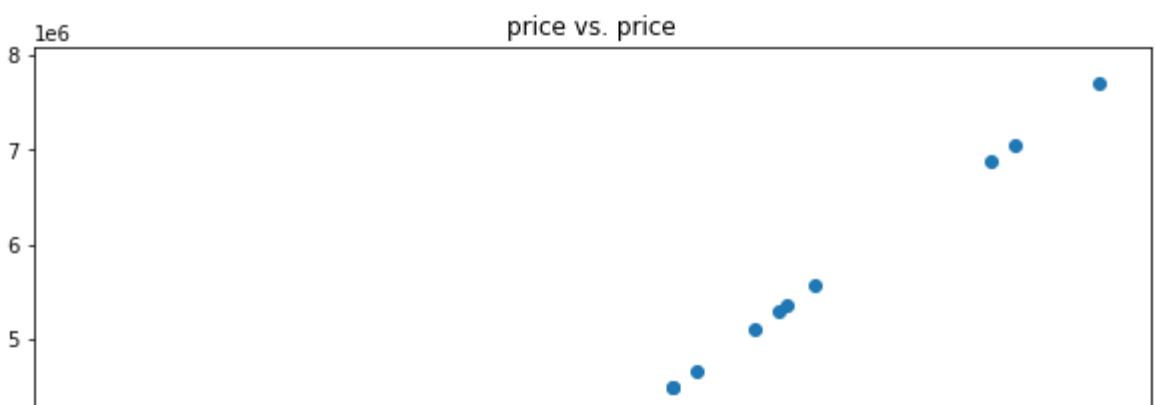
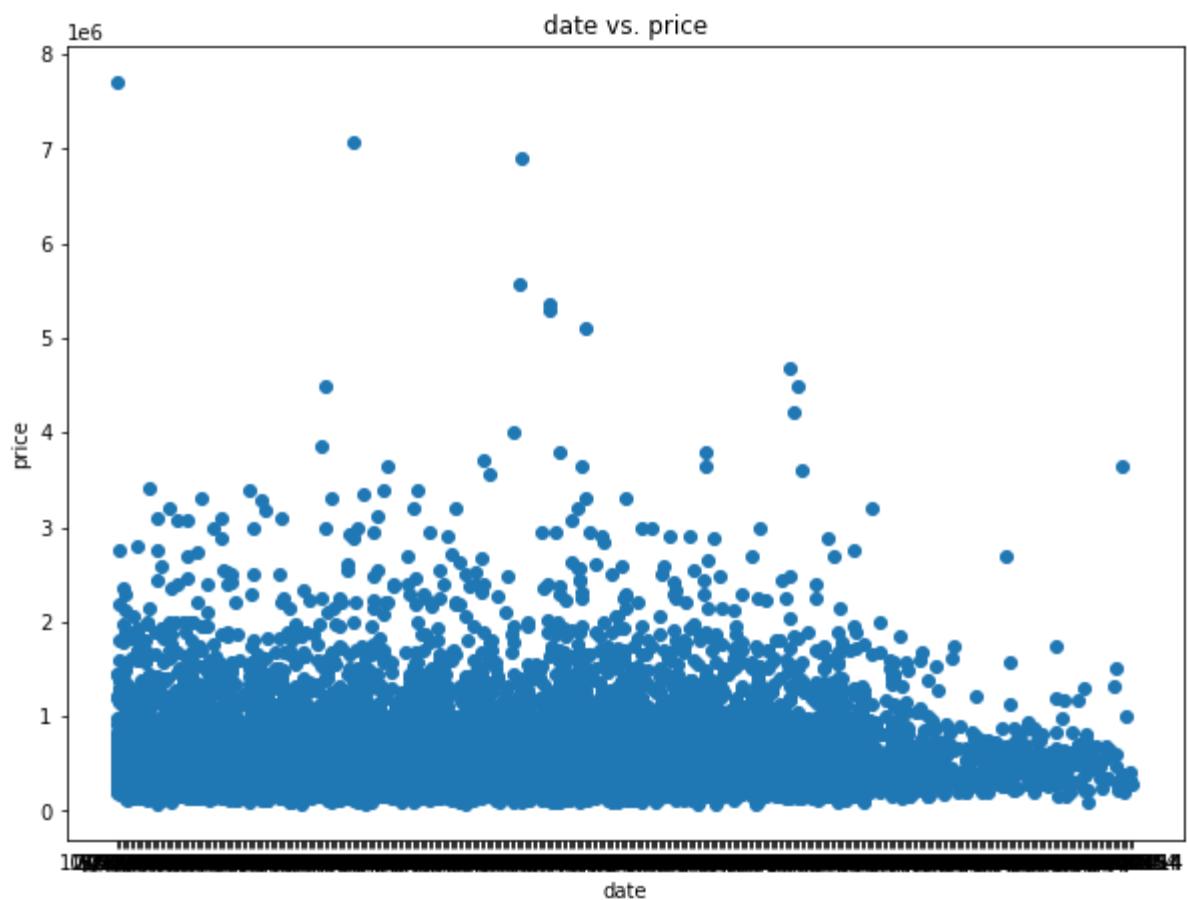
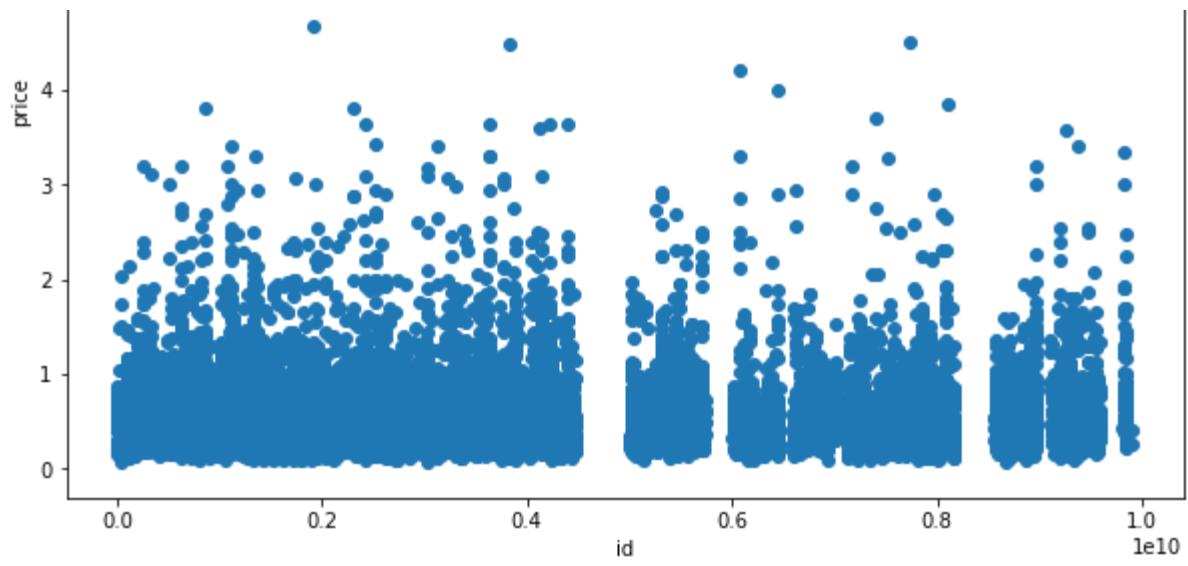
```
In [33]: import seaborn as sns

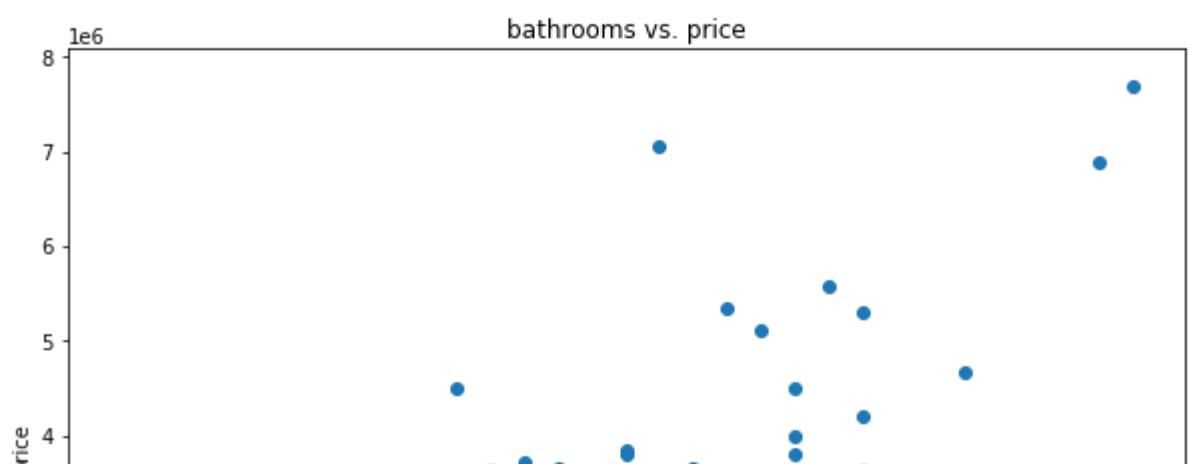
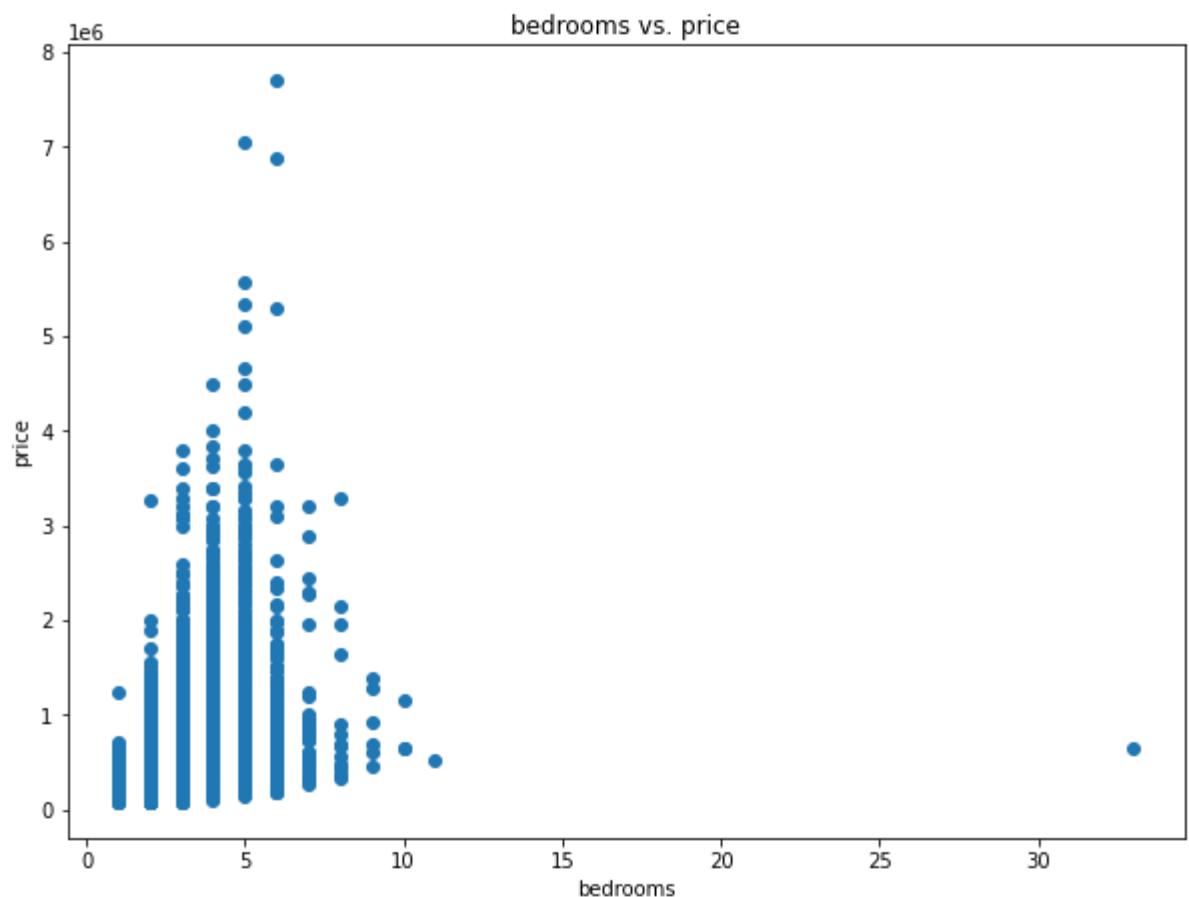
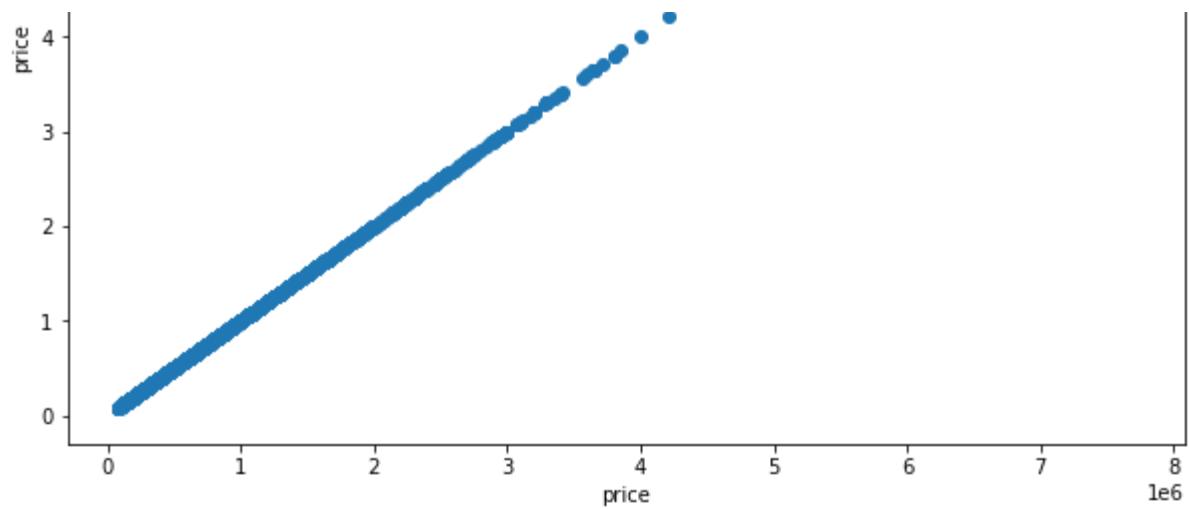
def plot(df, target='price'):
    fig, ax = plt.subplots(nrows = len(df.columns), figsize=(10,200))

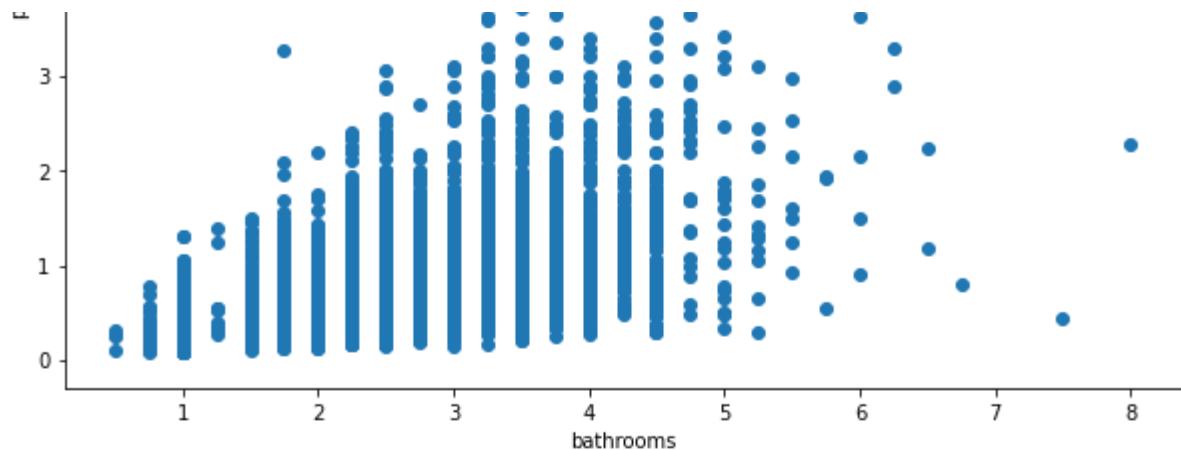
    for i, col in enumerate(df.columns):
        sns.lmplot(x=col, y=target, data=df)
        ax[i].scatter(df[col], df[target])
        ax[i].set_xlabel(col)
        ax[i].set_ylabel(target)
        ax[i].set_title(f'{col} vs. {target}'")
```

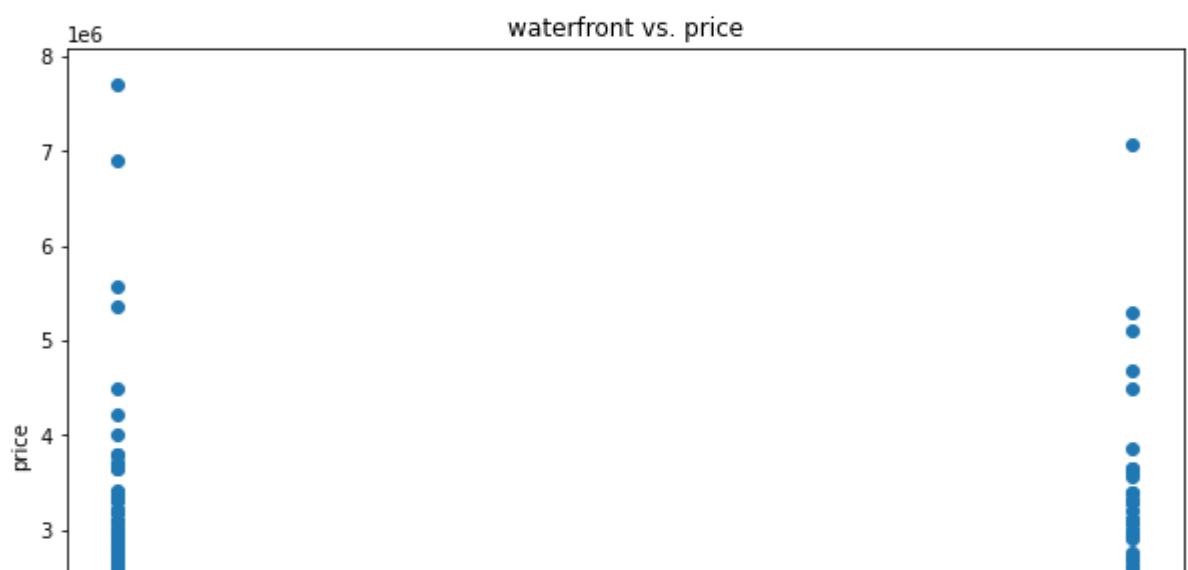
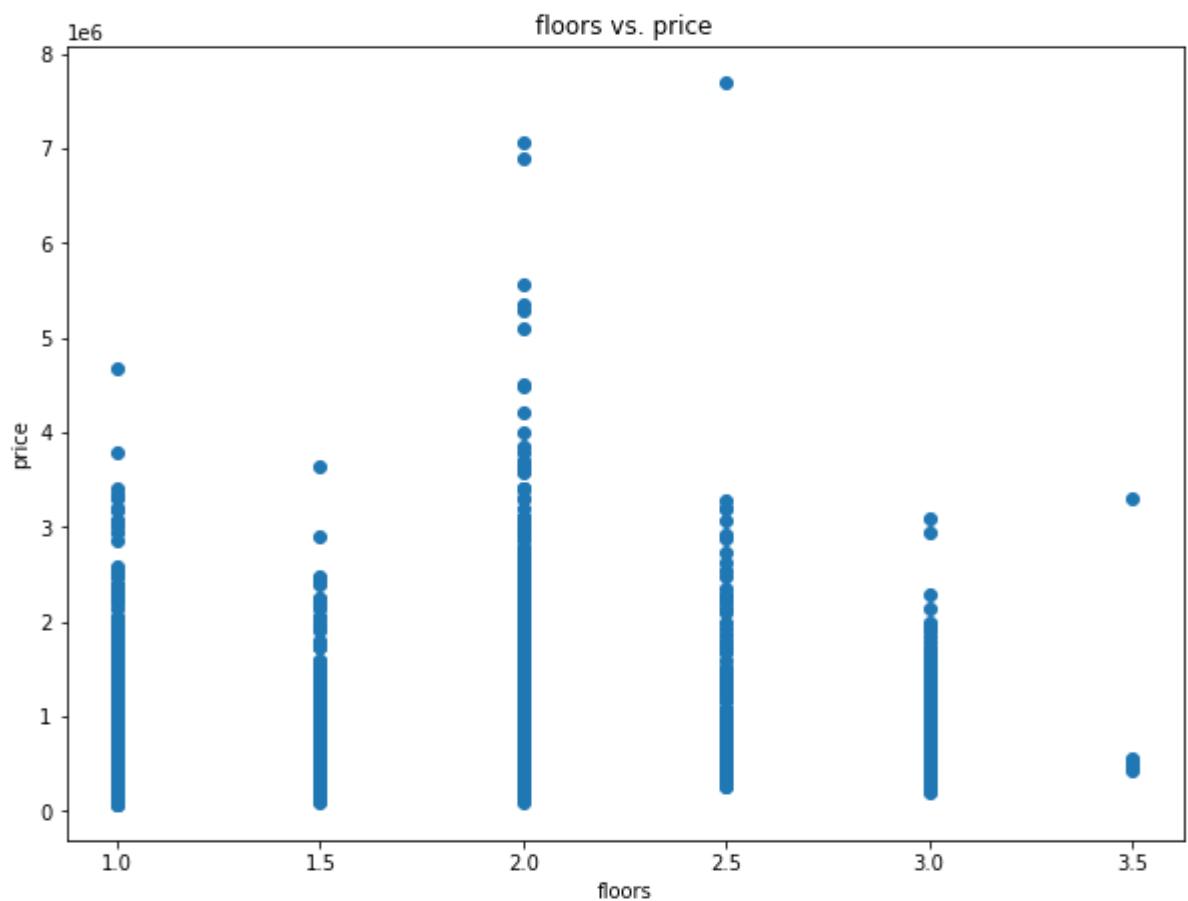
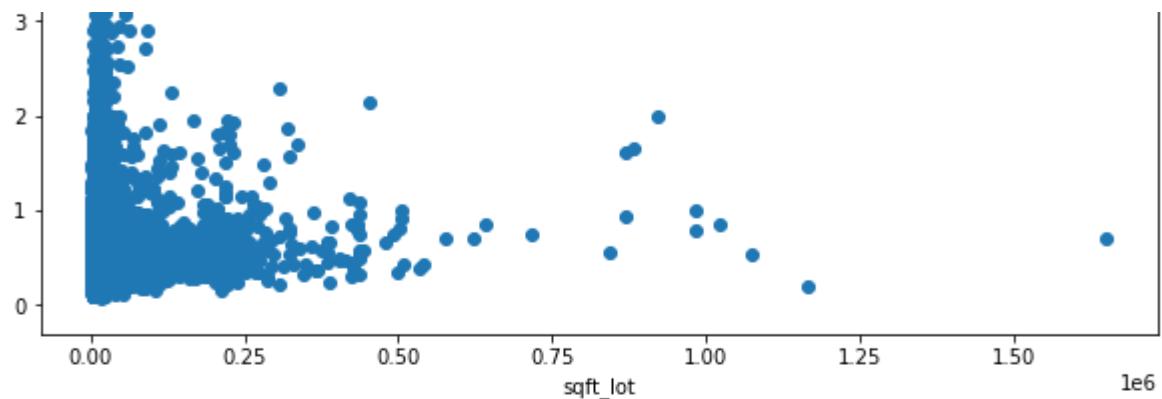
```
In [34]: plot(df=df, target='price')
```

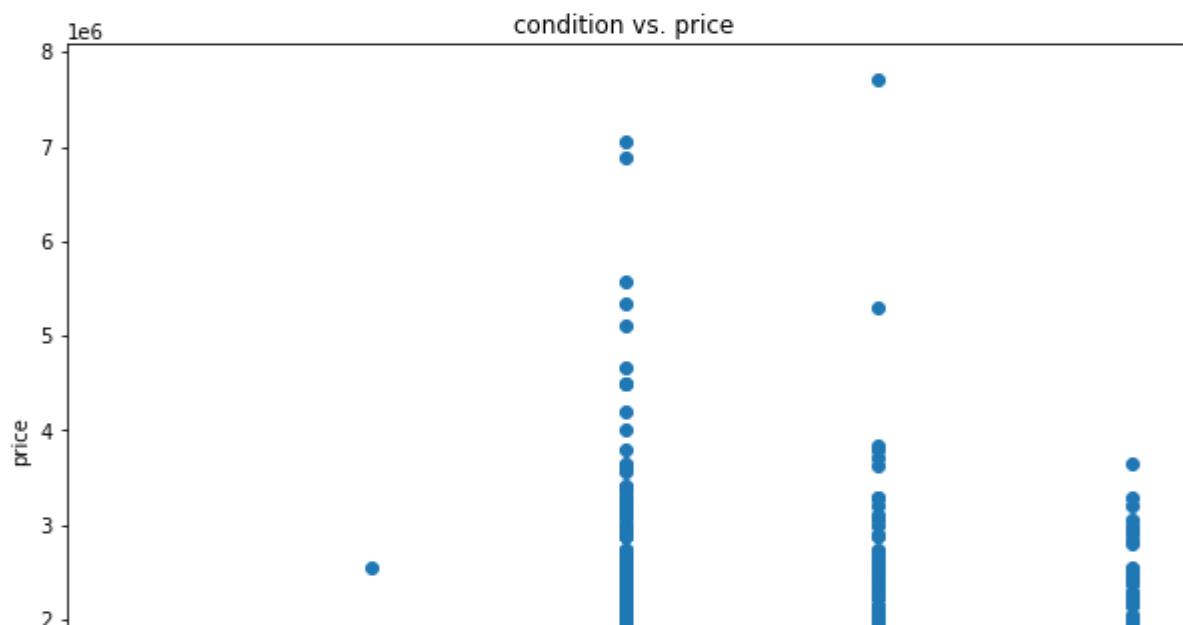
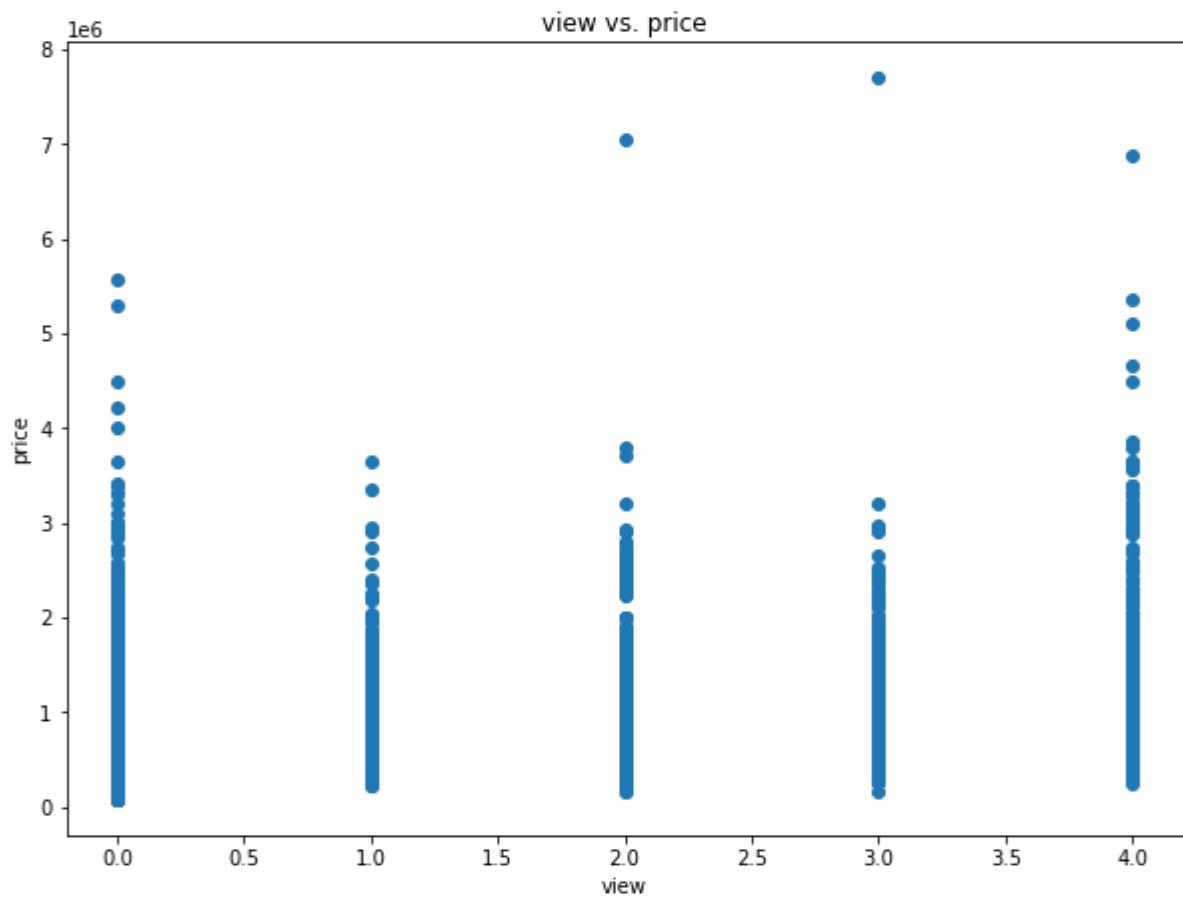
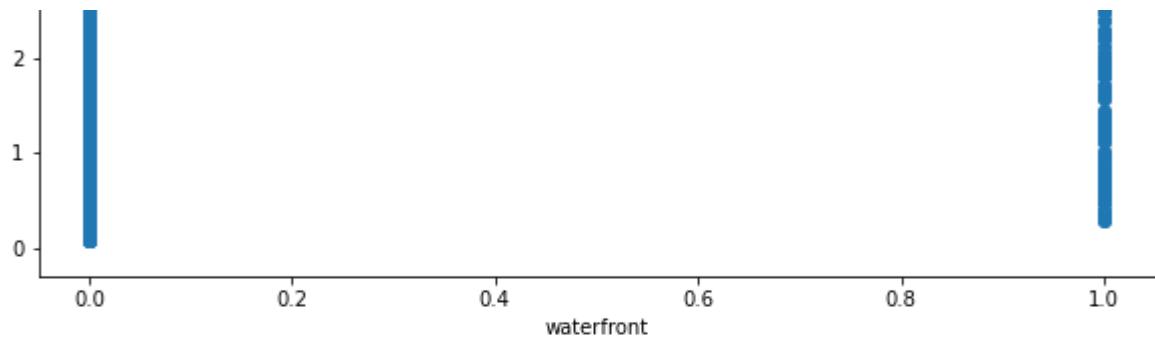


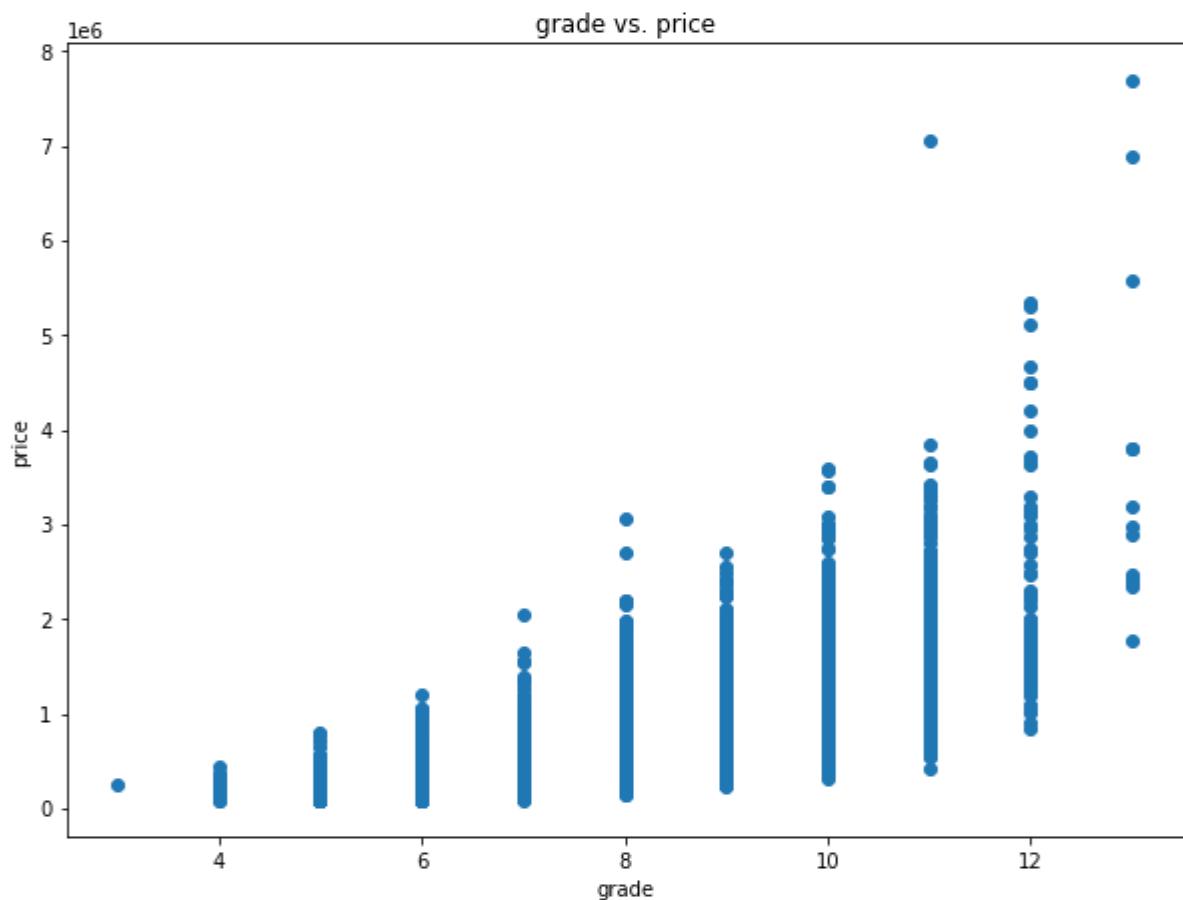
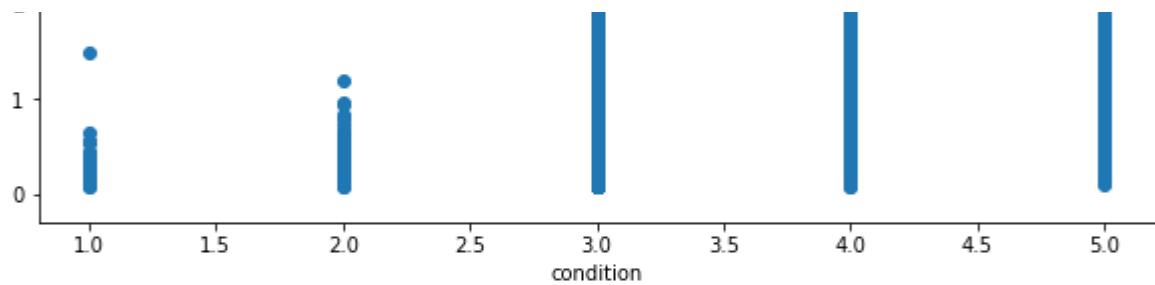


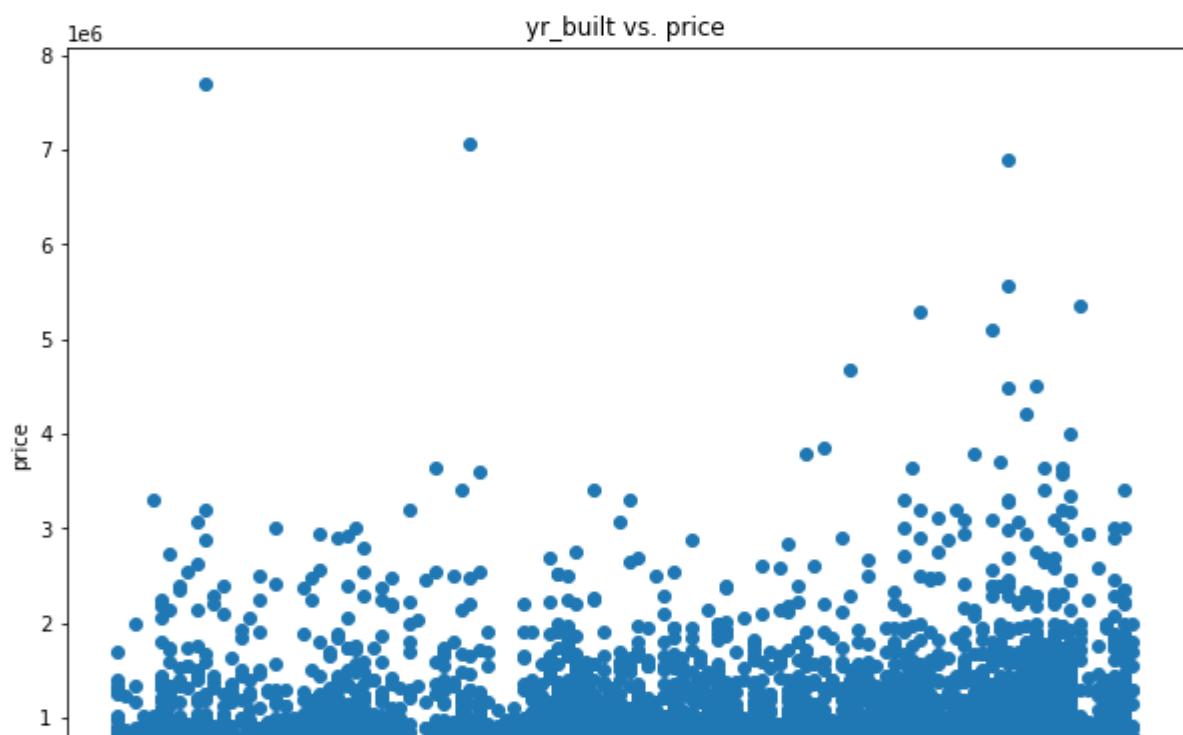
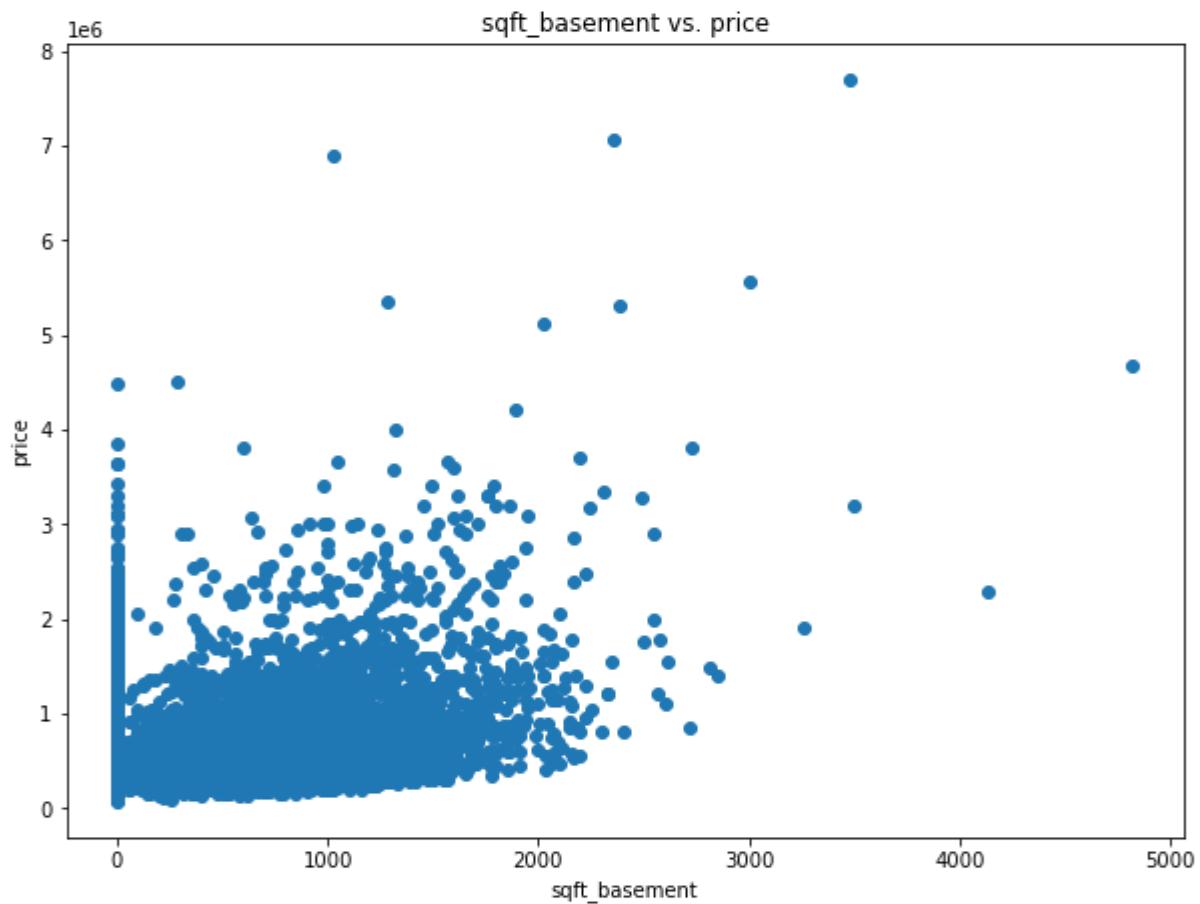
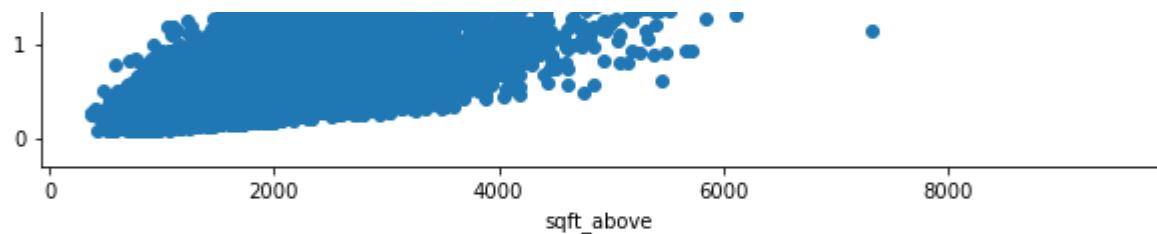


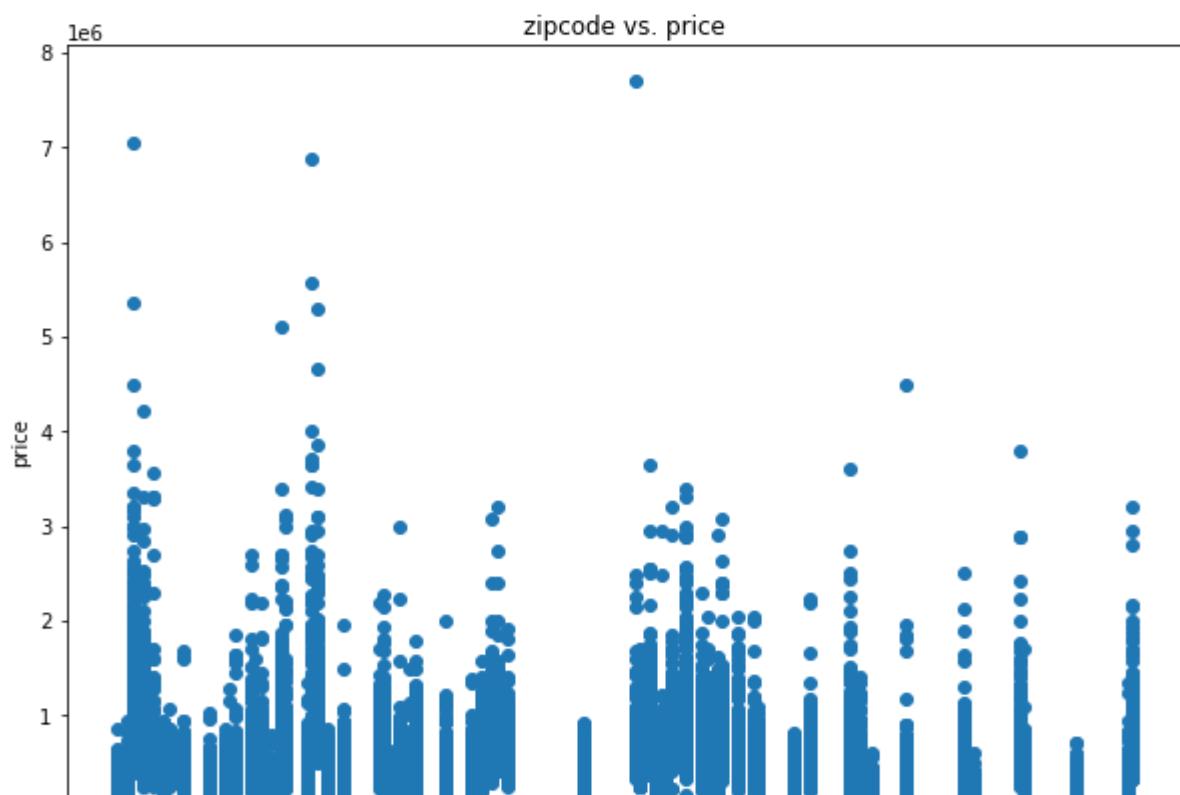
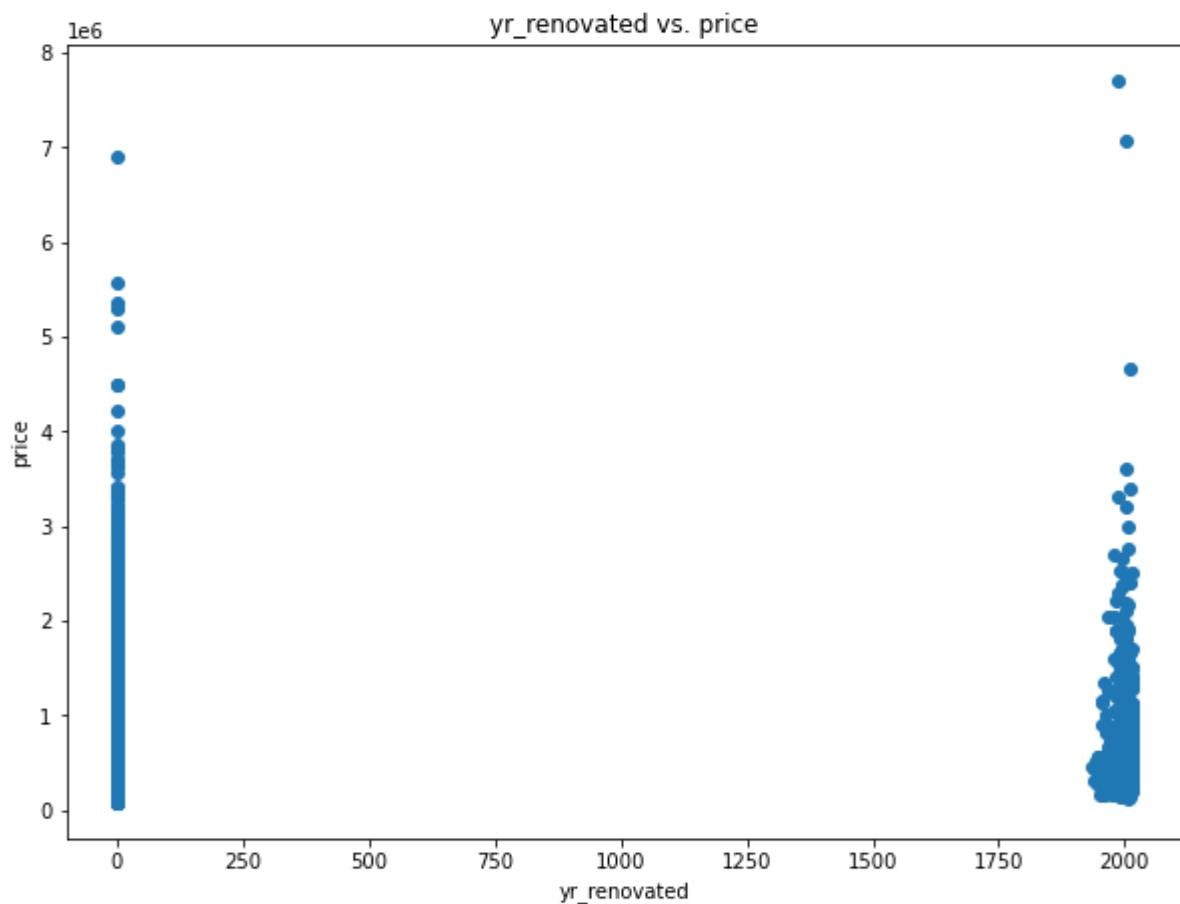
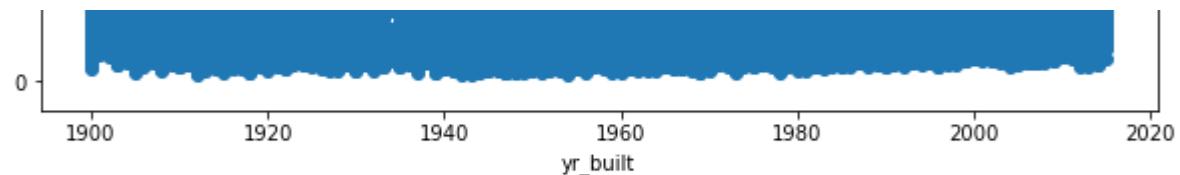


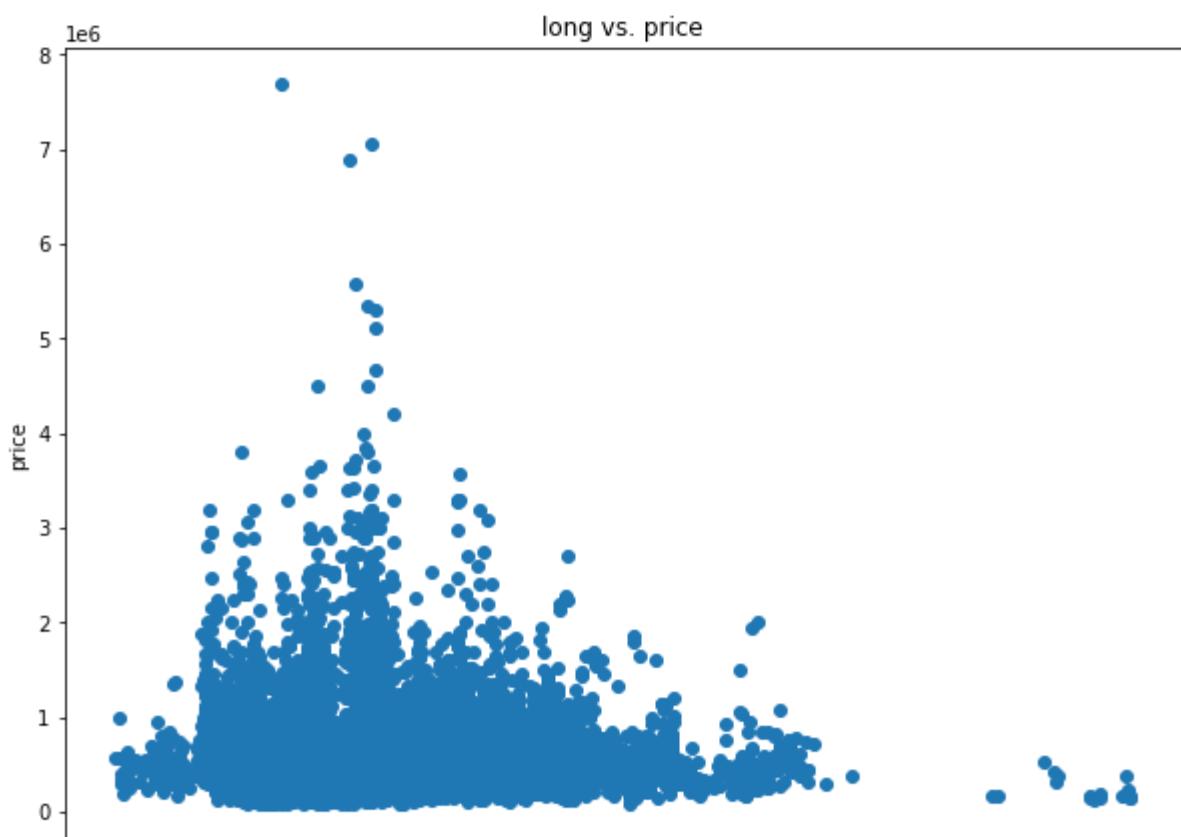
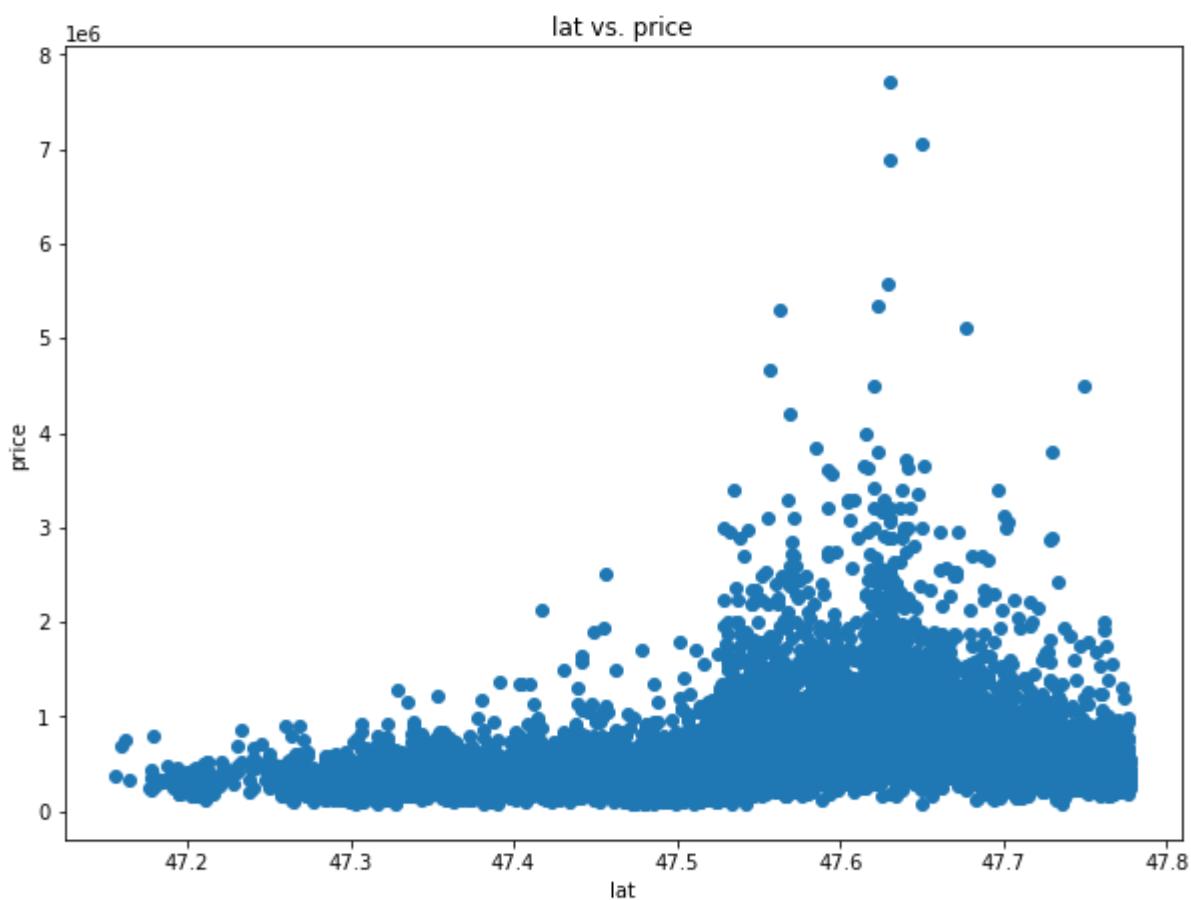
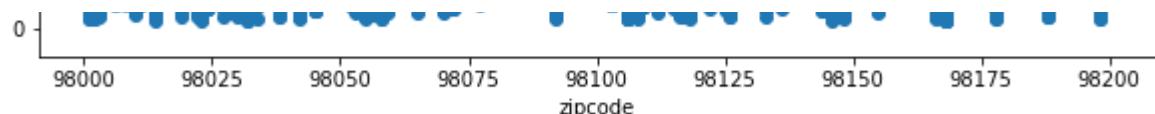


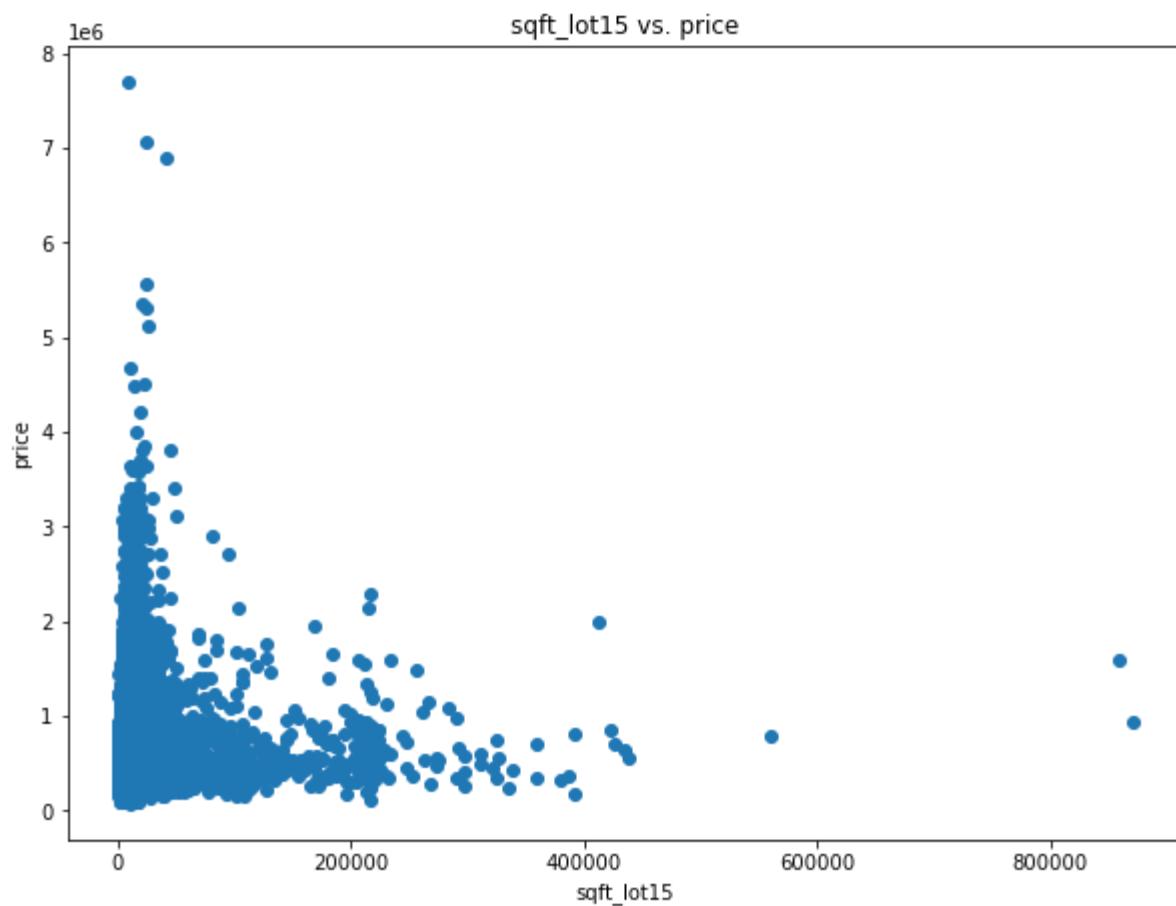
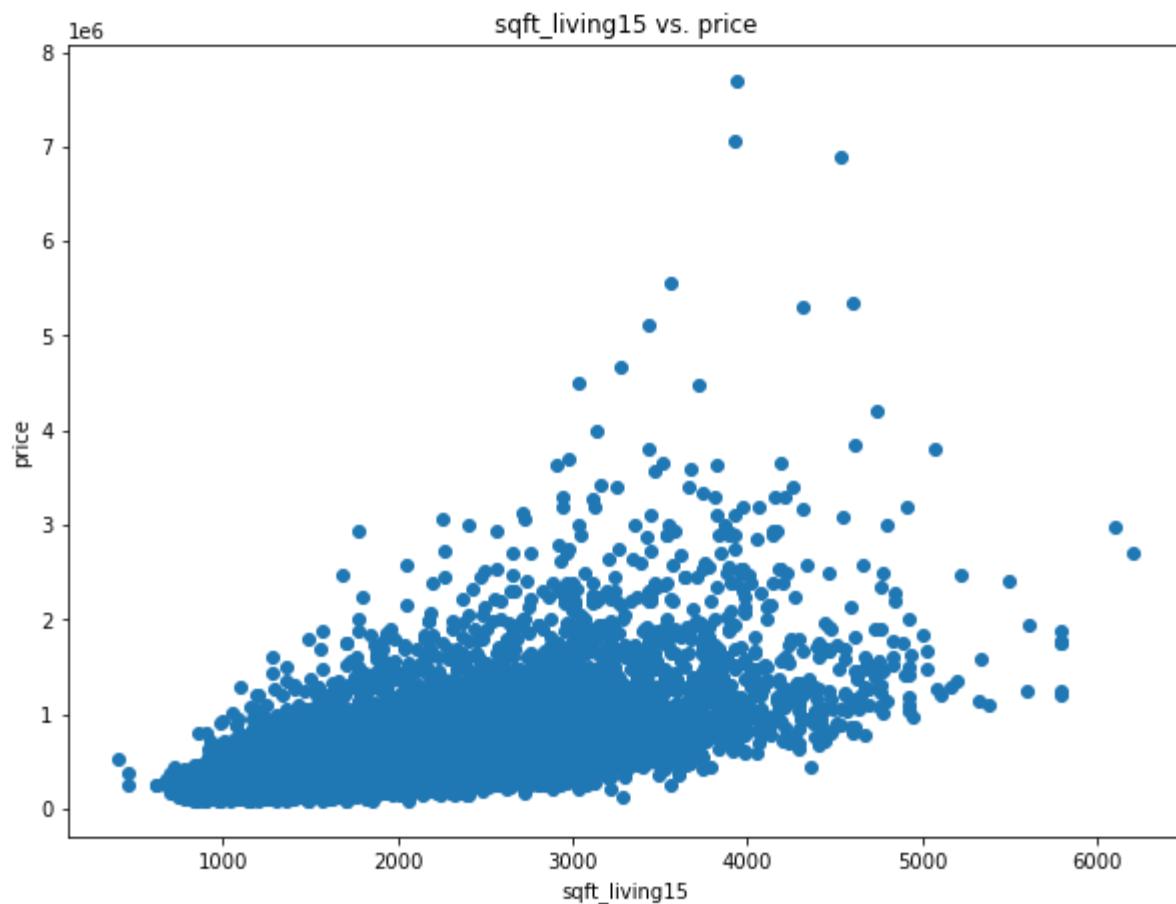


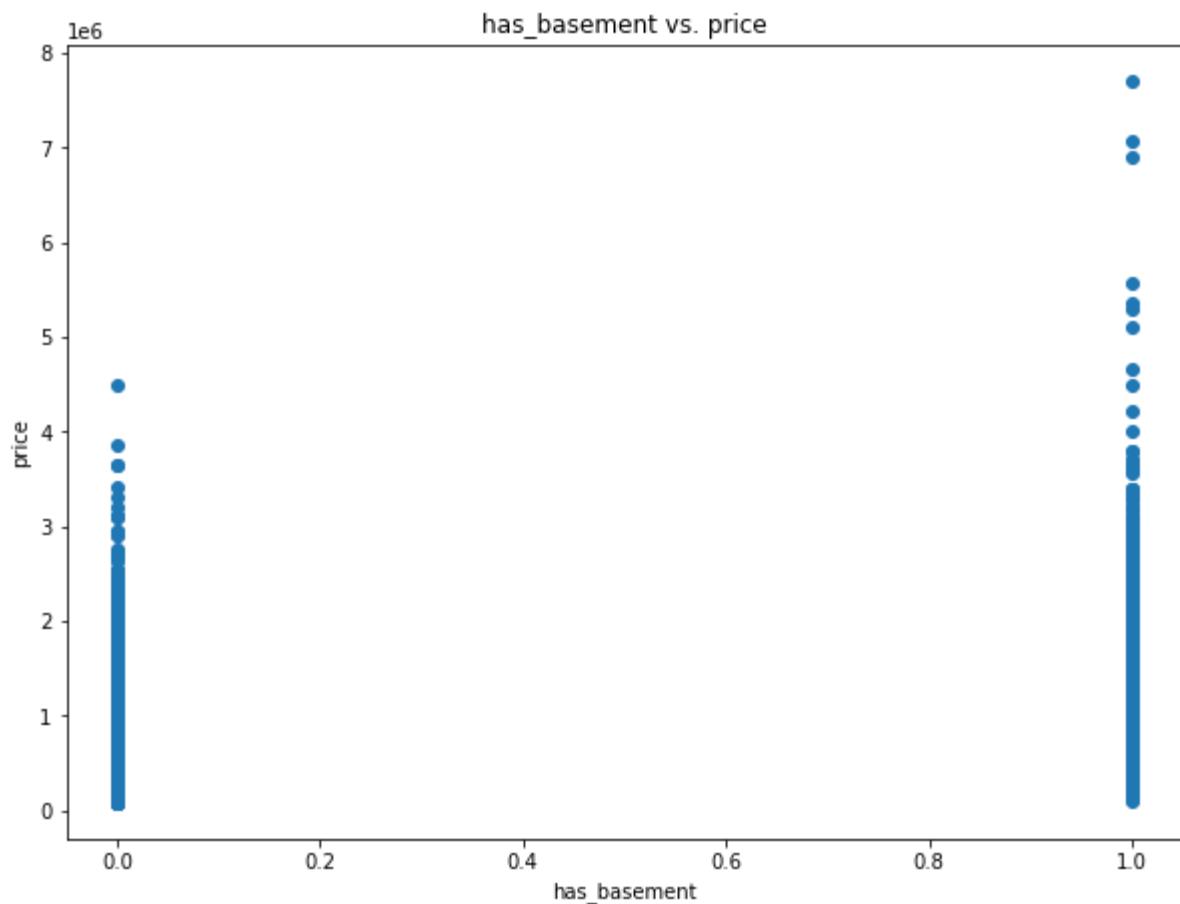
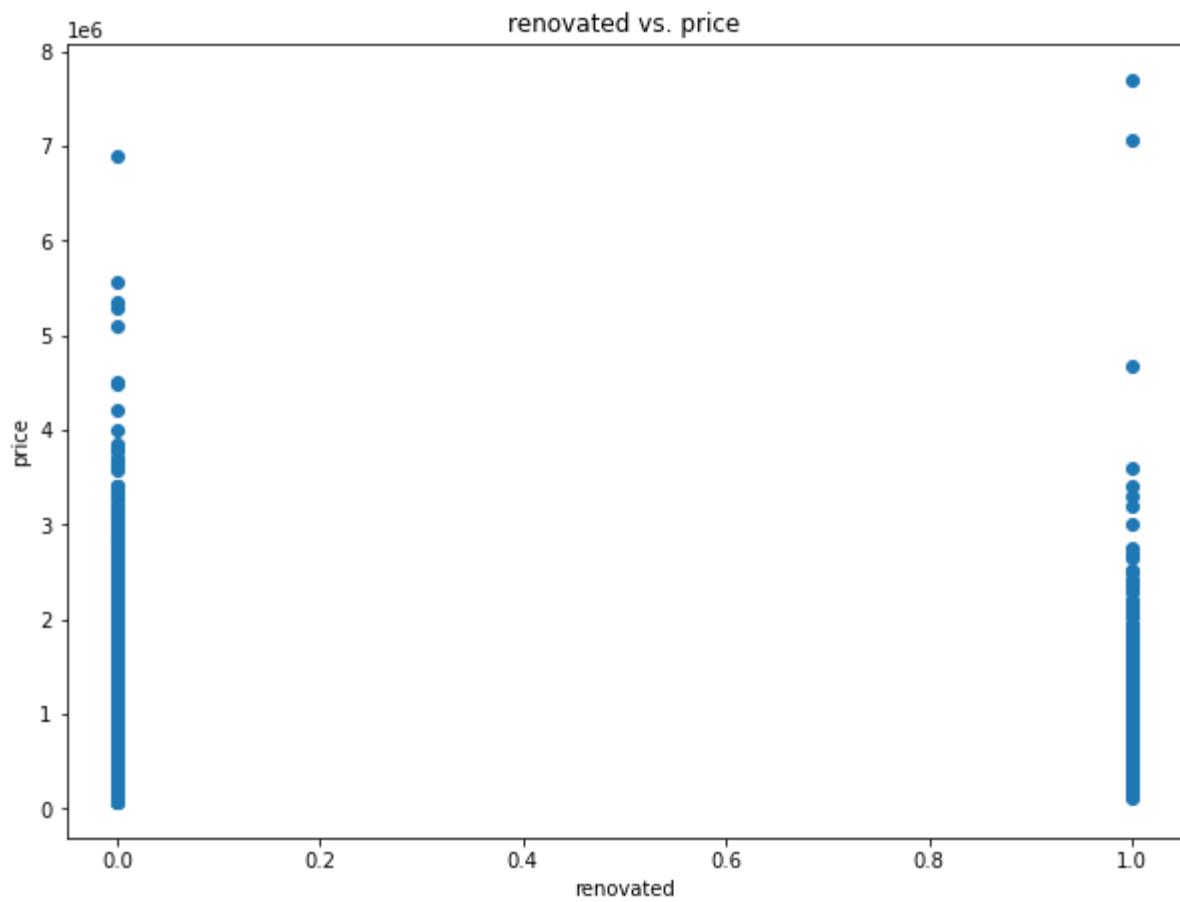








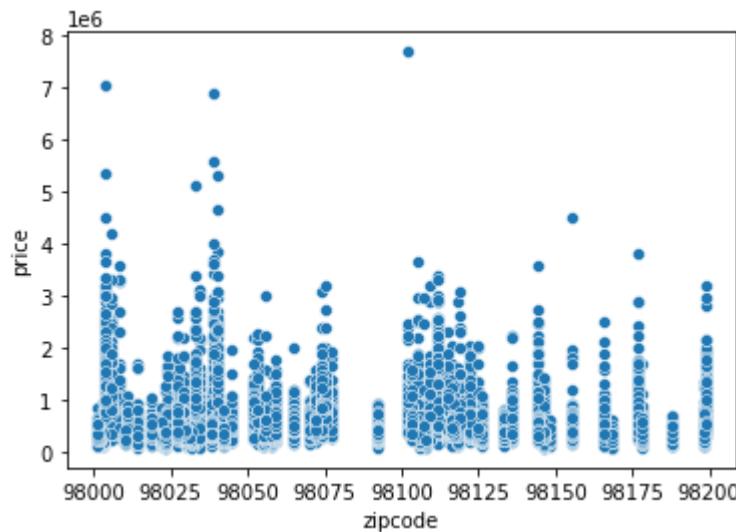




It seems as though there are certain columns that do **not** have a linear relationship with price. Based on the above scatter plots, at first glance, sqft\_lot, sqft\_lot15, bedrooms, floors, condition, and zipcode all have non-linear relationships with price. However, since we are to use an MLR model and would like to keep this information in to begin with, and then iteratively drop them if our models suggest high p-values, we will keep these columns as they are for the time being. The only exception is for categorical columns which we will need to one hot encode. Zipcode is clearly a categorical column that will need to be one hot encoded. It can be argued that columns such as grade, condition and bedrooms etc. can be treated as categorical columns but we can explore this later. Lastly, we can treat the binary columns such as has\_basement and renovated as they are already one hot encoded.

## One Hot Encoding - zipcode

```
In [35]: sns.scatterplot(x=df['zipcode'], y=df['price']);
```



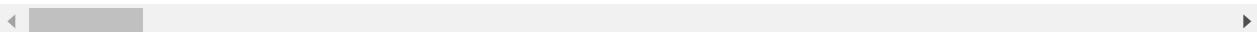
Zipcode is clearly a categorical column with no linear relationship with our target: 'price'. In order to keep this information in our model, we need to one hot encode this column.

```
In [36]: encoder = OneHotEncoder(sparse=False, drop='first')
cat_cols=['zipcode']
data_ohe = encoder.fit_transform(df[cat_cols])
df_ohe = pd.DataFrame(data_ohe, columns=encoder.get_feature_names(cat_cols), index=df.index)
df_ohe
```

| Out[36]:     | zipcode_98002 | zipcode_98003 | zipcode_98004 | zipcode_98005 | zipcode_98006 | zipcode_98007 | z   |
|--------------|---------------|---------------|---------------|---------------|---------------|---------------|-----|
| <b>0</b>     | 0.0           | 0.0           | 0.0           | 0.0           | 0.0           | 0.0           | 0.0 |
| <b>1</b>     | 0.0           | 0.0           | 0.0           | 0.0           | 0.0           | 0.0           | 0.0 |
| <b>2</b>     | 0.0           | 0.0           | 0.0           | 0.0           | 0.0           | 0.0           | 0.0 |
| <b>3</b>     | 0.0           | 0.0           | 0.0           | 0.0           | 0.0           | 0.0           | 0.0 |
| <b>4</b>     | 0.0           | 0.0           | 0.0           | 0.0           | 0.0           | 0.0           | 0.0 |
| <b>...</b>   | ...           | ...           | ...           | ...           | ...           | ...           | ... |
| <b>21592</b> | 0.0           | 0.0           | 0.0           | 0.0           | 0.0           | 0.0           | 0.0 |

|              | zipcode_98002 | zipcode_98003 | zipcode_98004 | zipcode_98005 | zipcode_98006 | zipcode_98007 | z   |
|--------------|---------------|---------------|---------------|---------------|---------------|---------------|-----|
| <b>21593</b> | 0.0           | 0.0           | 0.0           | 0.0           | 0.0           | 0.0           | 0.0 |
| <b>21594</b> | 0.0           | 0.0           | 0.0           | 0.0           | 0.0           | 0.0           | 0.0 |
| <b>21595</b> | 0.0           | 0.0           | 0.0           | 0.0           | 0.0           | 0.0           | 0.0 |
| <b>21596</b> | 0.0           | 0.0           | 0.0           | 0.0           | 0.0           | 0.0           | 0.0 |

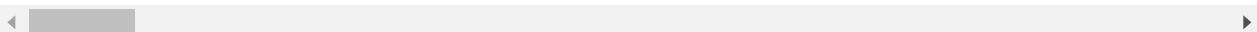
21597 rows × 69 columns



```
In [37]: df_ohe = pd.concat([df.drop('zipcode', axis=1), df_ohe], axis=1)
df_ohe
```

|              | id         | date       | price    | bedrooms | bathrooms | sqft_living | sqft_lot | floors | waterfront |
|--------------|------------|------------|----------|----------|-----------|-------------|----------|--------|------------|
| <b>0</b>     | 7129300520 | 10/13/2014 | 221900.0 | 3        | 1.00      | 1180        | 5650     | 1.0    | 0.         |
| <b>1</b>     | 6414100192 | 12/9/2014  | 538000.0 | 3        | 2.25      | 2570        | 7242     | 2.0    | 0.         |
| <b>2</b>     | 5631500400 | 2/25/2015  | 180000.0 | 2        | 1.00      | 770         | 10000    | 1.0    | 0.         |
| <b>3</b>     | 2487200875 | 12/9/2014  | 604000.0 | 4        | 3.00      | 1960        | 5000     | 1.0    | 0.         |
| <b>4</b>     | 1954400510 | 2/18/2015  | 510000.0 | 3        | 2.00      | 1680        | 8080     | 1.0    | 0.         |
| ...          | ...        | ...        | ...      | ...      | ...       | ...         | ...      | ...    | ...        |
| <b>21592</b> | 263000018  | 5/21/2014  | 360000.0 | 3        | 2.50      | 1530        | 1131     | 3.0    | 0.         |
| <b>21593</b> | 6600060120 | 2/23/2015  | 400000.0 | 4        | 2.50      | 2310        | 5813     | 2.0    | 0.         |
| <b>21594</b> | 1523300141 | 6/23/2014  | 402101.0 | 2        | 0.75      | 1020        | 1350     | 2.0    | 0.         |
| <b>21595</b> | 291310100  | 1/16/2015  | 400000.0 | 3        | 2.50      | 1600        | 2388     | 2.0    | 0.         |
| <b>21596</b> | 1523300157 | 10/15/2014 | 325000.0 | 2        | 0.75      | 1020        | 1076     | 2.0    | 0.         |

21597 rows × 91 columns



## Multicollinearity Check

### Correlation Matrices - Before

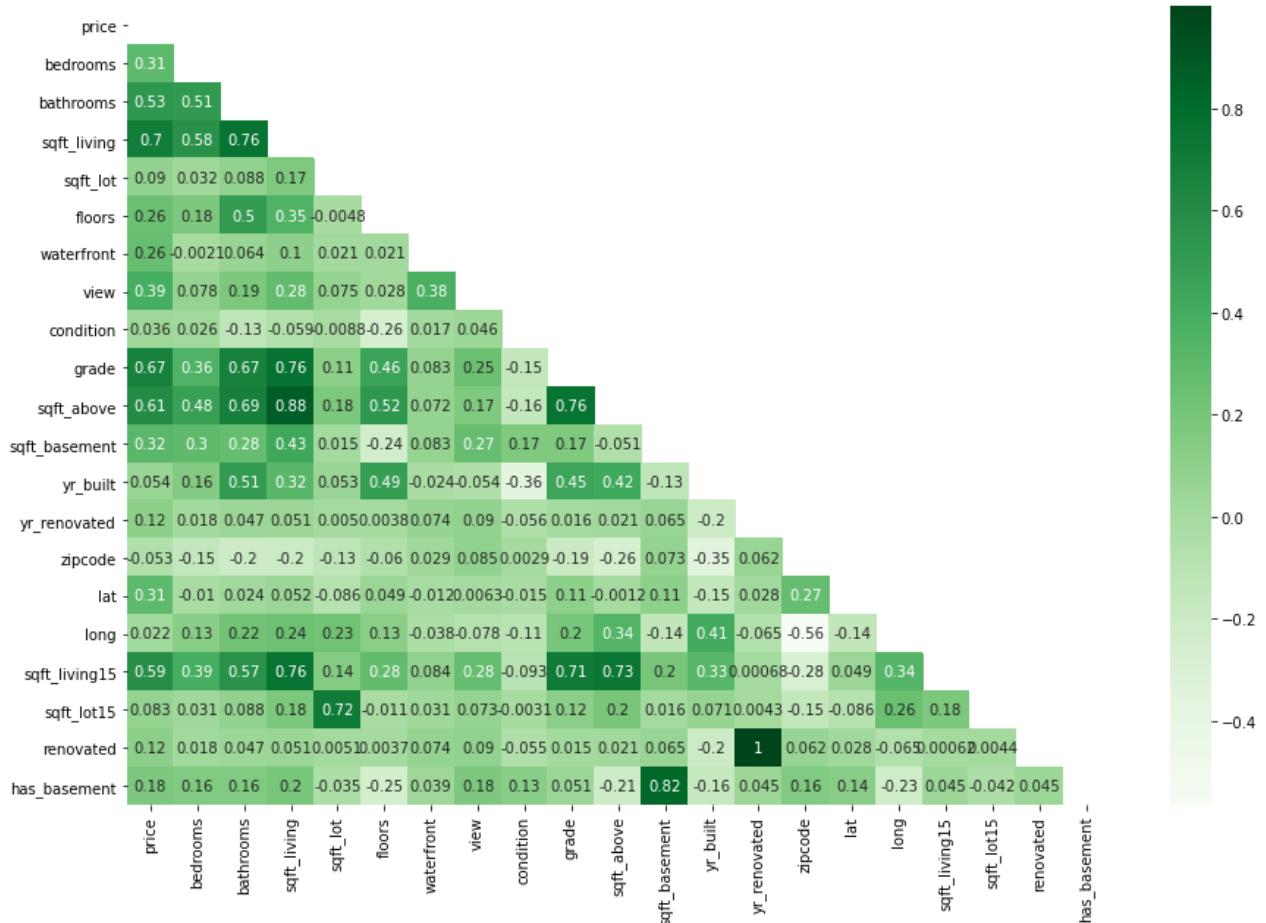
It looks like there is multicollinearity between the different columns. To address these I will be creating matrices and removing columns that may be causing the issues.

### With Target Column: price

```
In [38]: df.drop(['id', 'date'], axis=1, inplace=True)
```

```
In [39]: mask = np.zeros_like(df.corr())
mask[np.triu_indices_from(mask)] = True
fig, ax = plt.subplots(figsize=(15,10))
sns.heatmap(df.corr(), annot=True, mask=mask, cmap='Greens')
```

Out[39]: &lt;AxesSubplot:&gt;



In [40]: df.corr()['price'].abs().sort\_values(ascending=False)

```

Out[40]: price          1.000000
sqft_living      0.701917
grade            0.667951
sqft_above        0.605368
sqft_living15    0.585241
bathrooms         0.525906
view              0.393497
sqft_basement     0.321108
bedrooms          0.308787
lat                0.306692
waterfront         0.264306
floors             0.256804
has_basement       0.178264
yr_renovated      0.117855
renovated          0.117543
sqft_lot           0.089876
sqft_lot15         0.082845
yr_built            0.053953
zipcode            0.053402
condition          0.036056
long               0.022036
Name: price, dtype: float64

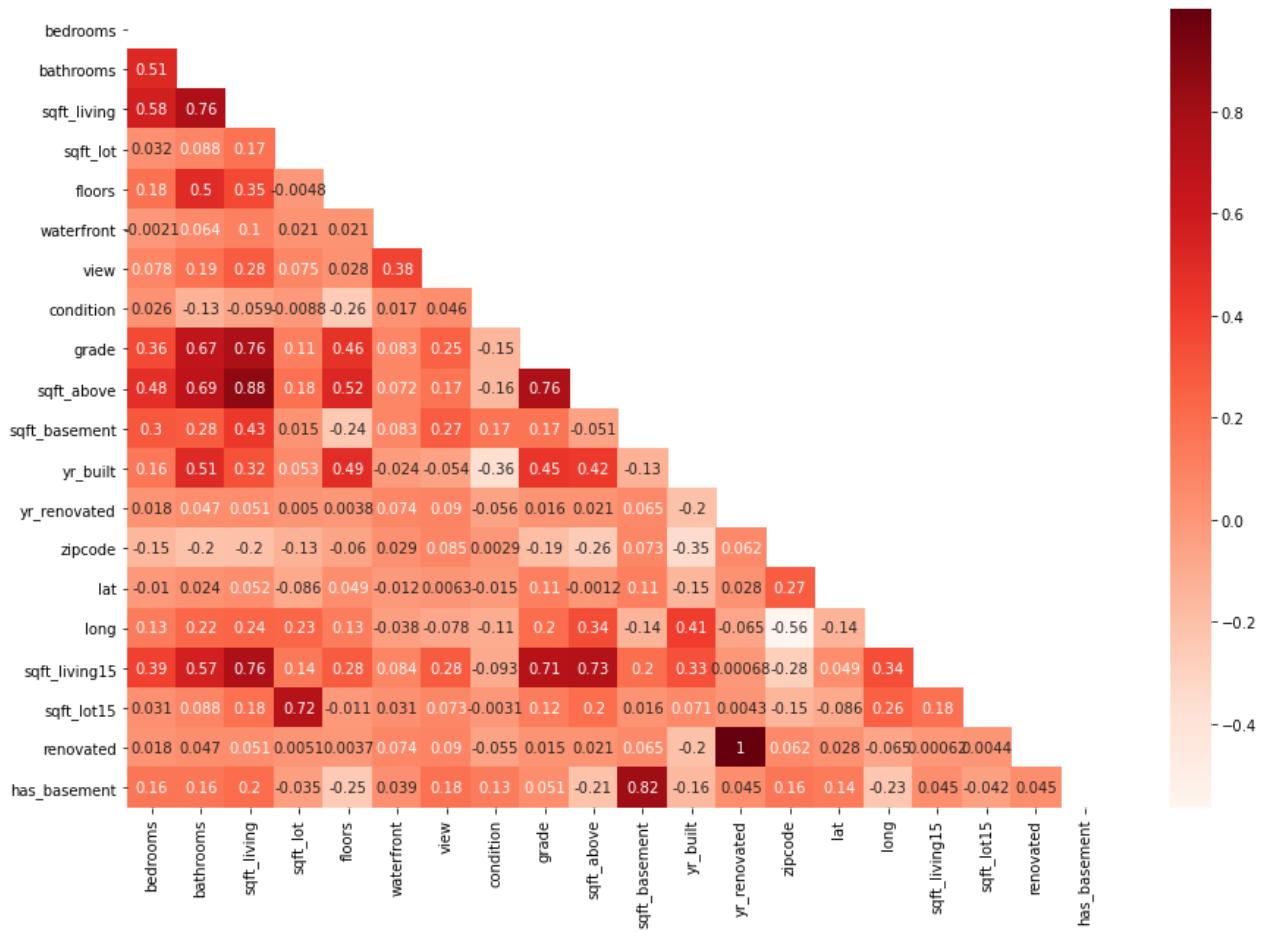
```

## Without Target Column: price

```
In [41]: mask = np.zeros_like(df.drop('price', axis=1).corr())
mask[np.triu_indices_from(mask)] = True
```

```
fig, ax = plt.subplots(figsize=(15,10))
sns.heatmap(df.drop('price', axis=1).corr(), annot=True, mask=mask, cmap='Reds')
```

Out[41]: &lt;AxesSubplot:&gt;



Looking at the correlation matrix of the original df (prior to being one hot encoded for simplicity of the visual), there seems to be strong correlations between a multitude of parameters when we take our cut-off point as 0.75. Dropping the 'sqft\_living' column will allow us to clear almost all the strong correlations. We are additionally not losing meaningful data since this column is equal to the sum of 'sqft\_above' and 'sqft\_basement'.

In [42]: `df.drop('sqft_living', axis=1, inplace=True)`

Since we have engineered a feature called 'has\_basement' using 'sqft\_basement' there is a high correlation between these columns as well. Since we are still keeping the basement information to a certain degree, we can also drop 'sqft\_basement'

In [43]: `df.drop('sqft_basement', axis=1, inplace=True)`

Similar to above, since we have engineered the column 'renovated' using 'yr\_renovated', we can also drop the 'yr\_renovated' column.

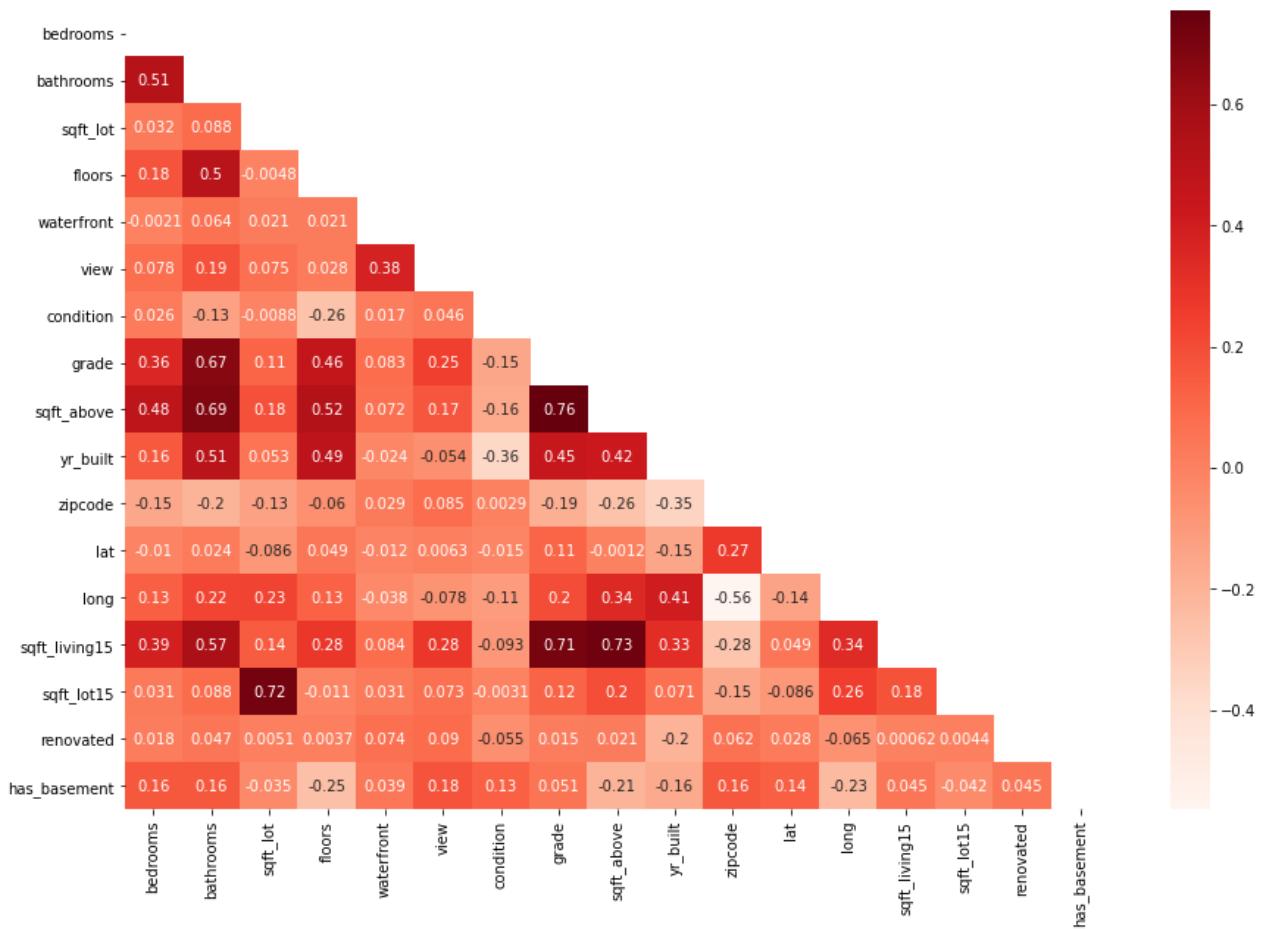
In [44]: `df.drop('yr_renovated', axis=1, inplace=True)`

## Correlation Matrix - After

In [45]: `mask = np.zeros_like(df.drop('price', axis=1).corr())`

```
mask[np.triu_indices_from(mask)] = True
fig, ax = plt.subplots(figsize=(15,10))
sns.heatmap(df.drop('price', axis=1).corr(), annot=True, mask=mask, cmap='Reds')
```

Out[45]: &lt;AxesSubplot:&gt;



Our correlation matrix looks much better with the exception of 'sqft\_above' and 'grade'. As this value is so close to our cut-off value of 0.75, we will be keeping it in to see if it causes any additional issues after adjustments to the model.

## MODEL

### Defining the Modelling Function - model\_lin\_reg

We are ready to fit the data with a linear regression model now. We chose this model because we were specifically asked to use this model. Since the building and improving on the model will be an iterative process, we can write a function that will not only create the model but also show us our QQ Plot as well as the residual information so we can check for normality and homoscedasticity.

In [46]:

```
def model_lin_reg(df, target='price'):

    features = ' + '.join(df.drop(target, axis=1).columns)
    f = f'{target}~{features}'
    model = smf.ols(f, df).fit()
    display(model.summary())
    fig, ax = plt.subplots(ncols=2, figsize=(15,5))
    sm.graphics.qqplot(model.resid, line='45', fit=True, ax=ax[0])
    sm.graphics.resid_plot(model.resid, line='45', fit=True, ax=ax[1])
```

```

sns.scatterplot(x=model.predict(df, transform=True), y=model.resid, ax=ax[1])
ax[1].set_ylabel('Residuals')
ax[1].set_xlabel('Predicted')
plt.axhline();
return model

```

```
In [47]: drop_cols = ['id', 'date', 'sqft_living', 'sqft_basement', 'yr_renovated']
df_ohe.drop(drop_cols, axis=1, inplace=True)
```

## Model #1 - Baseline

```
In [48]: model_lin_reg(df=df_ohe)
```

OLS Regression Results

|                          |                  |                            |             |
|--------------------------|------------------|----------------------------|-------------|
| <b>Dep. Variable:</b>    | price            | <b>R-squared:</b>          | 0.800       |
| <b>Model:</b>            | OLS              | <b>Adj. R-squared:</b>     | 0.799       |
| <b>Method:</b>           | Least Squares    | <b>F-statistic:</b>        | 1010.       |
| <b>Date:</b>             | Sat, 01 May 2021 | <b>Prob (F-statistic):</b> | 0.00        |
| <b>Time:</b>             | 16:50:20         | <b>Log-Likelihood:</b>     | -2.9003e+05 |
| <b>No. Observations:</b> | 21597            | <b>AIC:</b>                | 5.802e+05   |
| <b>Df Residuals:</b>     | 21511            | <b>BIC:</b>                | 5.809e+05   |
| <b>Df Model:</b>         | 85               |                            |             |
| <b>Covariance Type:</b>  | nonrobust        |                            |             |

|                      | <b>coef</b> | <b>std err</b> | <b>t</b> | <b>P&gt; t </b> | <b>[0.025</b> | <b>0.975]</b> |
|----------------------|-------------|----------------|----------|-----------------|---------------|---------------|
| <b>Intercept</b>     | -2.422e+07  | 6.3e+06        | -3.845   | 0.000           | -3.66e+07     | -1.19e+07     |
| <b>bedrooms</b>      | -1.763e+04  | 1544.455       | -11.413  | 0.000           | -2.07e+04     | -1.46e+04     |
| <b>bathrooms</b>     | 4.377e+04   | 2643.295       | 16.557   | 0.000           | 3.86e+04      | 4.89e+04      |
| <b>sqft_lot</b>      | 0.2560      | 0.039          | 6.499    | 0.000           | 0.179         | 0.333         |
| <b>floors</b>        | -6.367e+04  | 3189.871       | -19.961  | 0.000           | -6.99e+04     | -5.74e+04     |
| <b>waterfront</b>    | 6.937e+05   | 1.51e+04       | 46.012   | 0.000           | 6.64e+05      | 7.23e+05      |
| <b>view</b>          | 6.292e+04   | 1770.266       | 35.543   | 0.000           | 5.95e+04      | 6.64e+04      |
| <b>condition</b>     | 2.922e+04   | 1965.696       | 14.866   | 0.000           | 2.54e+04      | 3.31e+04      |
| <b>grade</b>         | 6.338e+04   | 1864.517       | 33.994   | 0.000           | 5.97e+04      | 6.7e+04       |
| <b>sqft_above</b>    | 190.3529    | 3.131          | 60.789   | 0.000           | 184.215       | 196.491       |
| <b>yr_built</b>      | -861.1302   | 65.764         | -13.094  | 0.000           | -990.033      | -732.228      |
| <b>lat</b>           | 2.032e+05   | 6.5e+04        | 3.125    | 0.002           | 7.58e+04      | 3.31e+05      |
| <b>long</b>          | -1.284e+05  | 4.67e+04       | -2.749   | 0.006           | -2.2e+05      | -3.69e+04     |
| <b>sqft_living15</b> | 21.3812     | 2.944          | 7.264    | 0.000           | 15.612        | 27.151        |
| <b>sqft_lot15</b>    | -0.0880     | 0.062          | -1.418   | 0.156           | -0.210        | 0.034         |

|                      |            |          |        |       |           |          |
|----------------------|------------|----------|--------|-------|-----------|----------|
| <b>renovated</b>     | 4.037e+04  | 6547.392 | 6.166  | 0.000 | 2.75e+04  | 5.32e+04 |
| <b>has_basement</b>  | 6.225e+04  | 3093.061 | 20.124 | 0.000 | 5.62e+04  | 6.83e+04 |
| <b>zipcode_98002</b> | 3.318e+04  | 1.48e+04 | 2.236  | 0.025 | 4093.329  | 6.23e+04 |
| <b>zipcode_98003</b> | -2.615e+04 | 1.33e+04 | -1.971 | 0.049 | -5.22e+04 | -149.702 |
| <b>zipcode_98004</b> | 7.234e+05  | 2.41e+04 | 30.009 | 0.000 | 6.76e+05  | 7.71e+05 |
| <b>zipcode_98005</b> | 2.487e+05  | 2.58e+04 | 9.653  | 0.000 | 1.98e+05  | 2.99e+05 |
| <b>zipcode_98006</b> | 2.355e+05  | 2.11e+04 | 11.178 | 0.000 | 1.94e+05  | 2.77e+05 |
| <b>zipcode_98007</b> | 1.937e+05  | 2.66e+04 | 7.286  | 0.000 | 1.42e+05  | 2.46e+05 |
| <b>zipcode_98008</b> | 2.025e+05  | 2.53e+04 | 8.018  | 0.000 | 1.53e+05  | 2.52e+05 |
| <b>zipcode_98010</b> | 9.554e+04  | 2.26e+04 | 4.223  | 0.000 | 5.12e+04  | 1.4e+05  |
| <b>zipcode_98011</b> | 3.791e+04  | 3.29e+04 | 1.154  | 0.249 | -2.65e+04 | 1.02e+05 |
| <b>zipcode_98014</b> | 8.408e+04  | 3.61e+04 | 2.330  | 0.020 | 1.33e+04  | 1.55e+05 |
| <b>zipcode_98019</b> | 4.459e+04  | 3.56e+04 | 1.253  | 0.210 | -2.52e+04 | 1.14e+05 |
| <b>zipcode_98022</b> | 3.069e+04  | 1.97e+04 | 1.562  | 0.118 | -7832.522 | 6.92e+04 |
| <b>zipcode_98023</b> | -4.989e+04 | 1.22e+04 | -4.088 | 0.000 | -7.38e+04 | -2.6e+04 |
| <b>zipcode_98024</b> | 1.431e+05  | 3.18e+04 | 4.504  | 0.000 | 8.08e+04  | 2.05e+05 |
| <b>zipcode_98027</b> | 1.547e+05  | 2.16e+04 | 7.153  | 0.000 | 1.12e+05  | 1.97e+05 |
| <b>zipcode_98028</b> | 3.521e+04  | 3.19e+04 | 1.103  | 0.270 | -2.73e+04 | 9.77e+04 |
| <b>zipcode_98029</b> | 1.881e+05  | 2.47e+04 | 7.614  | 0.000 | 1.4e+05   | 2.36e+05 |
| <b>zipcode_98030</b> | -1278.3007 | 1.46e+04 | -0.088 | 0.930 | -2.99e+04 | 2.73e+04 |
| <b>zipcode_98031</b> | 2733.1401  | 1.52e+04 | 0.180  | 0.857 | -2.7e+04  | 3.25e+04 |
| <b>zipcode_98032</b> | -6121.9357 | 1.76e+04 | -0.347 | 0.728 | -4.07e+04 | 2.84e+04 |
| <b>zipcode_98033</b> | 2.972e+05  | 2.74e+04 | 10.857 | 0.000 | 2.44e+05  | 3.51e+05 |
| <b>zipcode_98034</b> | 1.22e+05   | 2.93e+04 | 4.157  | 0.000 | 6.45e+04  | 1.8e+05  |
| <b>zipcode_98038</b> | 4.76e+04   | 1.64e+04 | 2.906  | 0.004 | 1.55e+04  | 7.97e+04 |
| <b>zipcode_98039</b> | 1.254e+06  | 3.26e+04 | 38.492 | 0.000 | 1.19e+06  | 1.32e+06 |
| <b>zipcode_98040</b> | 4.697e+05  | 2.13e+04 | 22.035 | 0.000 | 4.28e+05  | 5.12e+05 |
| <b>zipcode_98042</b> | 9795.9407  | 1.4e+04  | 0.702  | 0.483 | -1.76e+04 | 3.72e+04 |
| <b>zipcode_98045</b> | 1.186e+05  | 3.03e+04 | 3.920  | 0.000 | 5.93e+04  | 1.78e+05 |
| <b>zipcode_98052</b> | 1.633e+05  | 2.79e+04 | 5.844  | 0.000 | 1.09e+05  | 2.18e+05 |
| <b>zipcode_98053</b> | 1.436e+05  | 2.99e+04 | 4.794  | 0.000 | 8.49e+04  | 2.02e+05 |
| <b>zipcode_98055</b> | 2.336e+04  | 1.69e+04 | 1.381  | 0.167 | -9798.967 | 5.65e+04 |
| <b>zipcode_98056</b> | 6.749e+04  | 1.84e+04 | 3.672  | 0.000 | 3.15e+04  | 1.04e+05 |
| <b>zipcode_98058</b> | 1.286e+04  | 1.6e+04  | 0.805  | 0.421 | -1.85e+04 | 4.42e+04 |
| <b>zipcode_98059</b> | 5.588e+04  | 1.8e+04  | 3.099  | 0.002 | 2.05e+04  | 9.12e+04 |

|                      |            |          |        |       |           |           |
|----------------------|------------|----------|--------|-------|-----------|-----------|
| <b>zipcode_98065</b> | 8.765e+04  | 2.79e+04 | 3.142  | 0.002 | 3.3e+04   | 1.42e+05  |
| <b>zipcode_98070</b> | -5.806e+04 | 2.13e+04 | -2.730 | 0.006 | -9.97e+04 | -1.64e+04 |
| <b>zipcode_98072</b> | 7.39e+04   | 3.27e+04 | 2.261  | 0.024 | 9836.718  | 1.38e+05  |
| <b>zipcode_98074</b> | 1.277e+05  | 2.65e+04 | 4.825  | 0.000 | 7.58e+04  | 1.8e+05   |
| <b>zipcode_98075</b> | 1.295e+05  | 2.54e+04 | 5.089  | 0.000 | 7.96e+04  | 1.79e+05  |
| <b>zipcode_98077</b> | 4.529e+04  | 3.4e+04  | 1.332  | 0.183 | -2.14e+04 | 1.12e+05  |
| <b>zipcode_98092</b> | -2.867e+04 | 1.33e+04 | -2.161 | 0.031 | -5.47e+04 | -2667.678 |
| <b>zipcode_98102</b> | 4.422e+05  | 2.83e+04 | 15.646 | 0.000 | 3.87e+05  | 4.98e+05  |
| <b>zipcode_98103</b> | 2.628e+05  | 2.65e+04 | 9.927  | 0.000 | 2.11e+05  | 3.15e+05  |
| <b>zipcode_98105</b> | 3.882e+05  | 2.72e+04 | 14.279 | 0.000 | 3.35e+05  | 4.41e+05  |
| <b>zipcode_98106</b> | 9.546e+04  | 1.96e+04 | 4.866  | 0.000 | 5.7e+04   | 1.34e+05  |
| <b>zipcode_98107</b> | 2.629e+05  | 2.73e+04 | 9.634  | 0.000 | 2.09e+05  | 3.16e+05  |
| <b>zipcode_98108</b> | 8.08e+04   | 2.17e+04 | 3.731  | 0.000 | 3.84e+04  | 1.23e+05  |
| <b>zipcode_98109</b> | 4.174e+05  | 2.81e+04 | 14.844 | 0.000 | 3.62e+05  | 4.73e+05  |
| <b>zipcode_98112</b> | 5.539e+05  | 2.5e+04  | 22.191 | 0.000 | 5.05e+05  | 6.03e+05  |
| <b>zipcode_98115</b> | 2.56e+05   | 2.69e+04 | 9.511  | 0.000 | 2.03e+05  | 3.09e+05  |
| <b>zipcode_98116</b> | 2.243e+05  | 2.19e+04 | 10.240 | 0.000 | 1.81e+05  | 2.67e+05  |
| <b>zipcode_98117</b> | 2.321e+05  | 2.73e+04 | 8.515  | 0.000 | 1.79e+05  | 2.85e+05  |
| <b>zipcode_98118</b> | 1.313e+05  | 1.91e+04 | 6.862  | 0.000 | 9.38e+04  | 1.69e+05  |
| <b>zipcode_98119</b> | 3.943e+05  | 2.66e+04 | 14.841 | 0.000 | 3.42e+05  | 4.46e+05  |
| <b>zipcode_98122</b> | 2.704e+05  | 2.37e+04 | 11.410 | 0.000 | 2.24e+05  | 3.17e+05  |
| <b>zipcode_98125</b> | 1.21e+05   | 2.91e+04 | 4.164  | 0.000 | 6.4e+04   | 1.78e+05  |
| <b>zipcode_98126</b> | 1.369e+05  | 2.01e+04 | 6.804  | 0.000 | 9.74e+04  | 1.76e+05  |
| <b>zipcode_98133</b> | 7.66e+04   | 3e+04    | 2.553  | 0.011 | 1.78e+04  | 1.35e+05  |
| <b>zipcode_98136</b> | 1.855e+05  | 2.06e+04 | 8.994  | 0.000 | 1.45e+05  | 2.26e+05  |
| <b>zipcode_98144</b> | 2.254e+05  | 2.2e+04  | 10.231 | 0.000 | 1.82e+05  | 2.69e+05  |
| <b>zipcode_98146</b> | 6.678e+04  | 1.84e+04 | 3.628  | 0.000 | 3.07e+04  | 1.03e+05  |
| <b>zipcode_98148</b> | 3.46e+04   | 2.5e+04  | 1.382  | 0.167 | -1.45e+04 | 8.37e+04  |
| <b>zipcode_98155</b> | 5.55e+04   | 3.12e+04 | 1.779  | 0.075 | -5656.754 | 1.17e+05  |
| <b>zipcode_98166</b> | 1.449e+04  | 1.68e+04 | 0.860  | 0.390 | -1.85e+04 | 4.75e+04  |
| <b>zipcode_98168</b> | 4.58e+04   | 1.78e+04 | 2.573  | 0.010 | 1.09e+04  | 8.07e+04  |
| <b>zipcode_98177</b> | 1.168e+05  | 3.13e+04 | 3.727  | 0.000 | 5.54e+04  | 1.78e+05  |
| <b>zipcode_98178</b> | 1.105e+04  | 1.84e+04 | 0.601  | 0.548 | -2.5e+04  | 4.71e+04  |
| <b>zipcode_98188</b> | 1.379e+04  | 1.89e+04 | 0.731  | 0.465 | -2.32e+04 | 5.08e+04  |
| <b>zipcode_98198</b> | -2.15e+04  | 1.43e+04 | -1.504 | 0.133 | -4.95e+04 | 6528.999  |

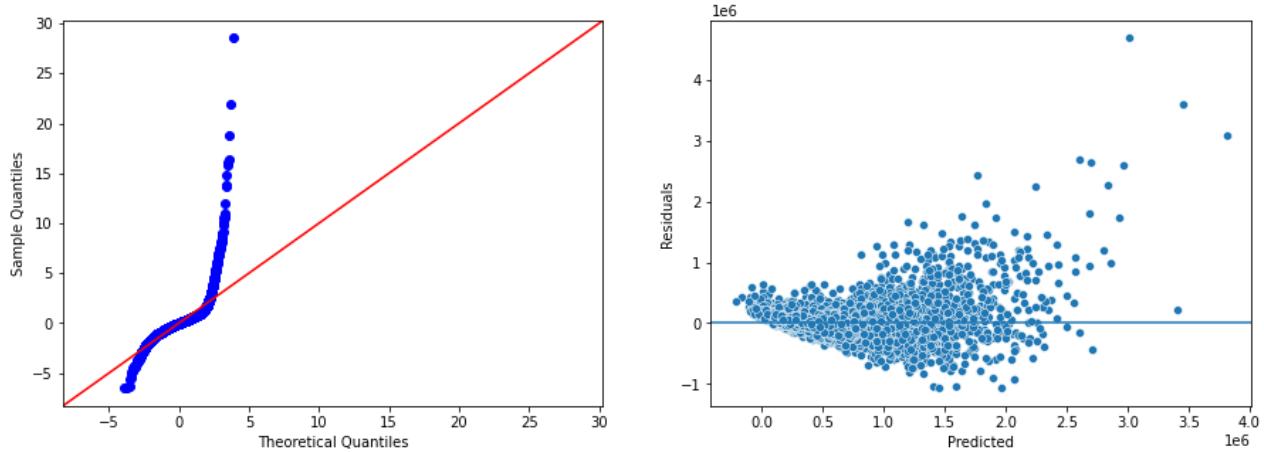
```
zipcode_98199 3.025e+05 2.59e+04 11.685 0.000 2.52e+05 3.53e+05
```

|                       |           |                          |             |
|-----------------------|-----------|--------------------------|-------------|
| <b>Omnibus:</b>       | 21922.754 | <b>Durbin-Watson:</b>    | 1.992       |
| <b>Prob(Omnibus):</b> | 0.000     | <b>Jarque-Bera (JB):</b> | 5315892.814 |
| <b>Skew:</b>          | 4.513     | <b>Prob(JB):</b>         | 0.00        |
| <b>Kurtosis:</b>      | 79.328    | <b>Cond. No.</b>         | 2.84e+08    |

Notes:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [2] The condition number is large, 2.84e+08. This might indicate that there are strong multicollinearity or other numerical problems.

```
Out[48]: <statsmodels.regression.linear_model.RegressionResultsWrapper at 0x24d046625b0>
```



Our baseline model has an R-squared value of 0.800 and an adjusted R-squared value of 0.799, which is fairly high. These will be two of various metrics we will be keeping an eye on to see how our model improves through iterations.

Another metric we will be looking for will be the p-values of our parameters (columns). If they are found insignificant (based on an alpha level of 0.05) we will be dropping these columns. In our baseline model, we are seeing that the sqft\_lot15 column has a higher p-value than 0.05 so its coefficient (its effect on our target) is not statistically significant and we can drop this column in its entirety.

The final metric are the residuals and how they fare with the MLR model's assumptions of normality and homoscedasticity. We are seeing here that they are definitely not normal according to our QQ plot and they are heteroscedastic. So, we are far from being done with our model.

As discussed above, sqft\_lot15 p-value is higher than our alpha of 0.05 meaning that this coefficient is not statistically significant, so we can go ahead and drop this parameter.

```
In [49]: df_ohe.drop('sqft_lot15', axis=1, inplace=True)
```

## Feature Engineering - has\_larger\_sqft\_than\_neighbors

We also wanted to create a feature to see if having a larger sqft (excluding the basement) would have any effect on the sale price of the home.

```
In [50]: df_ohe['has_larger_sqft_than_neighbors'] = df_ohe['sqft_living15'] < df_ohe['sqft_above']
df_ohe['has_larger_sqft_than_neighbors'] = df_ohe['has_larger_sqft_than_neighbors'].ast
df_ohe.drop('sqft_living15', axis=1, inplace=True)
df_ohe['has_larger_sqft_than_neighbors'].value_counts()
```

```
Out[50]: 0    15756
1     5841
Name: has_larger_sqft_than_neighbors, dtype: int64
```

Let's see how our model changed after we removed the statistically insignificant column 'sqft\_lot15' and added the new feature 'has\_larger\_sqft\_than\_neighbors'.

```
In [51]: model_lin_reg(df=df_ohe)
```

| OLS Regression Results |                  |                     |             |       |           |           |        |  |  |
|------------------------|------------------|---------------------|-------------|-------|-----------|-----------|--------|--|--|
| Dep. Variable:         | price            | R-squared:          | 0.800       |       |           |           |        |  |  |
| Model:                 | OLS              | Adj. R-squared:     | 0.800       |       |           |           |        |  |  |
| Method:                | Least Squares    | F-statistic:        | 1027.       |       |           |           |        |  |  |
| Date:                  | Sat, 01 May 2021 | Prob (F-statistic): | 0.00        |       |           |           |        |  |  |
| Time:                  | 16:50:21         | Log-Likelihood:     | -2.8999e+05 |       |           |           |        |  |  |
| No. Observations:      | 21597            | AIC:                | 5.801e+05   |       |           |           |        |  |  |
| Df Residuals:          | 21512            | BIC:                | 5.808e+05   |       |           |           |        |  |  |
| Df Model:              | 84               |                     |             |       |           |           |        |  |  |
| Covariance Type:       | nonrobust        |                     |             |       |           |           |        |  |  |
|                        |                  | coef                | std err     | t     | P> t      | [0.025    | 0.975] |  |  |
| <b>Intercept</b>       | -2.425e+07       | 6.28e+06            | -3.863      | 0.000 | -3.66e+07 | -1.19e+07 |        |  |  |
| <b>bedrooms</b>        | -1.68e+04        | 1541.928            | -10.899     | 0.000 | -1.98e+04 | -1.38e+04 |        |  |  |
| <b>bathrooms</b>       | 4.332e+04        | 2638.007            | 16.422      | 0.000 | 3.82e+04  | 4.85e+04  |        |  |  |
| <b>sqft_lot</b>        | 0.2141           | 0.030               | 7.166       | 0.000 | 0.156     | 0.273     |        |  |  |
| <b>floors</b>          | -6.274e+04       | 3179.284            | -19.734     | 0.000 | -6.9e+04  | -5.65e+04 |        |  |  |
| <b>waterfront</b>      | 6.911e+05        | 1.5e+04             | 45.966      | 0.000 | 6.62e+05  | 7.21e+05  |        |  |  |
| <b>view</b>            | 6.39e+04         | 1737.190            | 36.784      | 0.000 | 6.05e+04  | 6.73e+04  |        |  |  |
| <b>condition</b>       | 2.946e+04        | 1960.975            | 15.022      | 0.000 | 2.56e+04  | 3.33e+04  |        |  |  |
| <b>grade</b>           | 6.282e+04        | 1829.256            | 34.343      | 0.000 | 5.92e+04  | 6.64e+04  |        |  |  |
| <b>sqft_above</b>      | 213.2508         | 3.187               | 66.919      | 0.000 | 207.005   | 219.497   |        |  |  |
| <b>yr_built</b>        | -881.6538        | 65.484              | -13.464     | 0.000 | -1010.006 | -753.301  |        |  |  |
| <b>lat</b>             | 1.978e+05        | 6.49e+04            | 3.048       | 0.002 | 7.06e+04  | 3.25e+05  |        |  |  |
| <b>long</b>            | -1.312e+05       | 4.65e+04            | -2.819      | 0.005 | -2.22e+05 | -4e+04    |        |  |  |
| <b>renovated</b>       | 4.136e+04        | 6531.960            | 6.332       | 0.000 | 2.86e+04  | 5.42e+04  |        |  |  |

|                      |            |          |        |       |           |           |
|----------------------|------------|----------|--------|-------|-----------|-----------|
| <b>has_basement</b>  | 6.363e+04  | 3051.663 | 20.850 | 0.000 | 5.76e+04  | 6.96e+04  |
| <b>zipcode_98002</b> | 3.298e+04  | 1.48e+04 | 2.228  | 0.026 | 3963.897  | 6.2e+04   |
| <b>zipcode_98003</b> | -2.736e+04 | 1.32e+04 | -2.067 | 0.039 | -5.33e+04 | -1414.683 |
| <b>zipcode_98004</b> | 7.254e+05  | 2.4e+04  | 30.163 | 0.000 | 6.78e+05  | 7.72e+05  |
| <b>zipcode_98005</b> | 2.515e+05  | 2.57e+04 | 9.786  | 0.000 | 2.01e+05  | 3.02e+05  |
| <b>zipcode_98006</b> | 2.386e+05  | 2.1e+04  | 11.363 | 0.000 | 1.97e+05  | 2.8e+05   |
| <b>zipcode_98007</b> | 1.964e+05  | 2.65e+04 | 7.405  | 0.000 | 1.44e+05  | 2.48e+05  |
| <b>zipcode_98008</b> | 2.039e+05  | 2.52e+04 | 8.092  | 0.000 | 1.55e+05  | 2.53e+05  |
| <b>zipcode_98010</b> | 9.662e+04  | 2.26e+04 | 4.281  | 0.000 | 5.24e+04  | 1.41e+05  |
| <b>zipcode_98011</b> | 4.068e+04  | 3.28e+04 | 1.241  | 0.215 | -2.36e+04 | 1.05e+05  |
| <b>zipcode_98014</b> | 8.781e+04  | 3.6e+04  | 2.438  | 0.015 | 1.72e+04  | 1.58e+05  |
| <b>zipcode_98019</b> | 4.763e+04  | 3.55e+04 | 1.341  | 0.180 | -2.2e+04  | 1.17e+05  |
| <b>zipcode_98022</b> | 3.037e+04  | 1.96e+04 | 1.548  | 0.122 | -8073.985 | 6.88e+04  |
| <b>zipcode_98023</b> | -5.138e+04 | 1.22e+04 | -4.219 | 0.000 | -7.53e+04 | -2.75e+04 |
| <b>zipcode_98024</b> | 1.409e+05  | 3.17e+04 | 4.448  | 0.000 | 7.88e+04  | 2.03e+05  |
| <b>zipcode_98027</b> | 1.551e+05  | 2.16e+04 | 7.190  | 0.000 | 1.13e+05  | 1.97e+05  |
| <b>zipcode_98028</b> | 3.818e+04  | 3.18e+04 | 1.199  | 0.231 | -2.42e+04 | 1.01e+05  |
| <b>zipcode_98029</b> | 1.889e+05  | 2.46e+04 | 7.667  | 0.000 | 1.41e+05  | 2.37e+05  |
| <b>zipcode_98030</b> | 366.4616   | 1.46e+04 | 0.025  | 0.980 | -2.82e+04 | 2.89e+04  |
| <b>zipcode_98031</b> | 1079.4560  | 1.52e+04 | 0.071  | 0.943 | -2.86e+04 | 3.08e+04  |
| <b>zipcode_98032</b> | -7082.3756 | 1.76e+04 | -0.403 | 0.687 | -4.16e+04 | 2.74e+04  |
| <b>zipcode_98033</b> | 2.993e+05  | 2.73e+04 | 10.957 | 0.000 | 2.46e+05  | 3.53e+05  |
| <b>zipcode_98034</b> | 1.23e+05   | 2.93e+04 | 4.199  | 0.000 | 6.56e+04  | 1.8e+05   |
| <b>zipcode_98038</b> | 4.753e+04  | 1.63e+04 | 2.909  | 0.004 | 1.55e+04  | 7.95e+04  |
| <b>zipcode_98039</b> | 1.254e+06  | 3.25e+04 | 38.575 | 0.000 | 1.19e+06  | 1.32e+06  |
| <b>zipcode_98040</b> | 4.749e+05  | 2.13e+04 | 22.348 | 0.000 | 4.33e+05  | 5.17e+05  |
| <b>zipcode_98042</b> | 9389.1280  | 1.39e+04 | 0.674  | 0.500 | -1.79e+04 | 3.67e+04  |
| <b>zipcode_98045</b> | 1.176e+05  | 3.02e+04 | 3.896  | 0.000 | 5.84e+04  | 1.77e+05  |
| <b>zipcode_98052</b> | 1.65e+05   | 2.79e+04 | 5.919  | 0.000 | 1.1e+05   | 2.2e+05   |
| <b>zipcode_98053</b> | 1.455e+05  | 2.99e+04 | 4.869  | 0.000 | 8.69e+04  | 2.04e+05  |
| <b>zipcode_98055</b> | 2.435e+04  | 1.69e+04 | 1.443  | 0.149 | -8737.813 | 5.74e+04  |
| <b>zipcode_98056</b> | 6.823e+04  | 1.83e+04 | 3.720  | 0.000 | 3.23e+04  | 1.04e+05  |
| <b>zipcode_98058</b> | 1.332e+04  | 1.59e+04 | 0.835  | 0.404 | -1.79e+04 | 4.46e+04  |
| <b>zipcode_98059</b> | 5.721e+04  | 1.8e+04  | 3.182  | 0.001 | 2.2e+04   | 9.25e+04  |
| <b>zipcode_98065</b> | 8.945e+04  | 2.78e+04 | 3.217  | 0.001 | 3.5e+04   | 1.44e+05  |

|                      |            |          |        |       |           |           |
|----------------------|------------|----------|--------|-------|-----------|-----------|
| <b>zipcode_98070</b> | -5.994e+04 | 2.11e+04 | -2.847 | 0.004 | -1.01e+05 | -1.87e+04 |
| <b>zipcode_98072</b> | 7.646e+04  | 3.26e+04 | 2.345  | 0.019 | 1.25e+04  | 1.4e+05   |
| <b>zipcode_98074</b> | 1.308e+05  | 2.64e+04 | 4.956  | 0.000 | 7.91e+04  | 1.83e+05  |
| <b>zipcode_98075</b> | 1.329e+05  | 2.54e+04 | 5.243  | 0.000 | 8.32e+04  | 1.83e+05  |
| <b>zipcode_98077</b> | 4.759e+04  | 3.39e+04 | 1.403  | 0.161 | -1.89e+04 | 1.14e+05  |
| <b>zipcode_98092</b> | -2.821e+04 | 1.32e+04 | -2.133 | 0.033 | -5.41e+04 | -2287.050 |
| <b>zipcode_98102</b> | 4.419e+05  | 2.82e+04 | 15.669 | 0.000 | 3.87e+05  | 4.97e+05  |
| <b>zipcode_98103</b> | 2.63e+05   | 2.64e+04 | 9.959  | 0.000 | 2.11e+05  | 3.15e+05  |
| <b>zipcode_98105</b> | 3.893e+05  | 2.71e+04 | 14.349 | 0.000 | 3.36e+05  | 4.42e+05  |
| <b>zipcode_98106</b> | 9.449e+04  | 1.96e+04 | 4.831  | 0.000 | 5.62e+04  | 1.33e+05  |
| <b>zipcode_98107</b> | 2.624e+05  | 2.72e+04 | 9.638  | 0.000 | 2.09e+05  | 3.16e+05  |
| <b>zipcode_98108</b> | 8.161e+04  | 2.16e+04 | 3.777  | 0.000 | 3.93e+04  | 1.24e+05  |
| <b>zipcode_98109</b> | 4.186e+05  | 2.81e+04 | 14.918 | 0.000 | 3.64e+05  | 4.74e+05  |
| <b>zipcode_98112</b> | 5.554e+05  | 2.49e+04 | 22.297 | 0.000 | 5.07e+05  | 6.04e+05  |
| <b>zipcode_98115</b> | 2.565e+05  | 2.69e+04 | 9.552  | 0.000 | 2.04e+05  | 3.09e+05  |
| <b>zipcode_98116</b> | 2.235e+05  | 2.18e+04 | 10.229 | 0.000 | 1.81e+05  | 2.66e+05  |
| <b>zipcode_98117</b> | 2.317e+05  | 2.72e+04 | 8.525  | 0.000 | 1.78e+05  | 2.85e+05  |
| <b>zipcode_98118</b> | 1.315e+05  | 1.91e+04 | 6.890  | 0.000 | 9.41e+04  | 1.69e+05  |
| <b>zipcode_98119</b> | 3.945e+05  | 2.65e+04 | 14.881 | 0.000 | 3.42e+05  | 4.46e+05  |
| <b>zipcode_98122</b> | 2.713e+05  | 2.36e+04 | 11.473 | 0.000 | 2.25e+05  | 3.18e+05  |
| <b>zipcode_98125</b> | 1.229e+05  | 2.9e+04  | 4.237  | 0.000 | 6.6e+04   | 1.8e+05   |
| <b>zipcode_98126</b> | 1.36e+05   | 2.01e+04 | 6.782  | 0.000 | 9.67e+04  | 1.75e+05  |
| <b>zipcode_98133</b> | 7.711e+04  | 2.99e+04 | 2.576  | 0.010 | 1.84e+04  | 1.36e+05  |
| <b>zipcode_98136</b> | 1.841e+05  | 2.06e+04 | 8.949  | 0.000 | 1.44e+05  | 2.24e+05  |
| <b>zipcode_98144</b> | 2.246e+05  | 2.2e+04  | 10.218 | 0.000 | 1.82e+05  | 2.68e+05  |
| <b>zipcode_98146</b> | 6.647e+04  | 1.84e+04 | 3.621  | 0.000 | 3.05e+04  | 1.02e+05  |
| <b>zipcode_98148</b> | 3.604e+04  | 2.5e+04  | 1.443  | 0.149 | -1.29e+04 | 8.5e+04   |
| <b>zipcode_98155</b> | 5.694e+04  | 3.11e+04 | 1.829  | 0.067 | -4090.436 | 1.18e+05  |
| <b>zipcode_98166</b> | 1.463e+04  | 1.68e+04 | 0.870  | 0.384 | -1.83e+04 | 4.76e+04  |
| <b>zipcode_98168</b> | 4.695e+04  | 1.78e+04 | 2.644  | 0.008 | 1.21e+04  | 8.18e+04  |
| <b>zipcode_98177</b> | 1.174e+05  | 3.13e+04 | 3.757  | 0.000 | 5.62e+04  | 1.79e+05  |
| <b>zipcode_98178</b> | 1.199e+04  | 1.83e+04 | 0.654  | 0.513 | -2.4e+04  | 4.8e+04   |
| <b>zipcode_98188</b> | 1.322e+04  | 1.88e+04 | 0.702  | 0.482 | -2.37e+04 | 5.01e+04  |
| <b>zipcode_98198</b> | -2.128e+04 | 1.43e+04 | -1.491 | 0.136 | -4.92e+04 | 6692.339  |
| <b>zipcode_98199</b> | 3.036e+05  | 2.58e+04 | 11.754 | 0.000 | 2.53e+05  | 3.54e+05  |

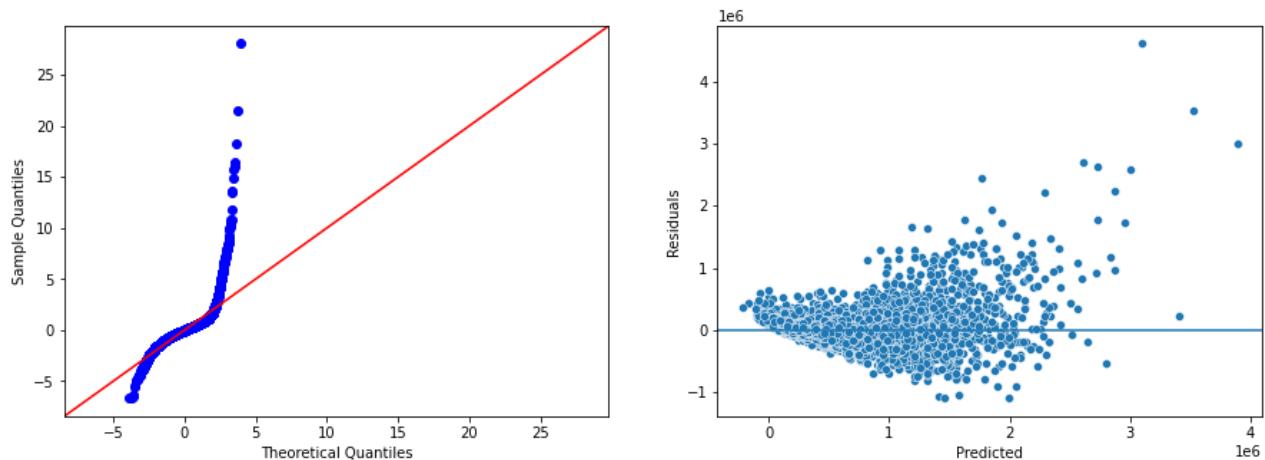
```
has_larger_sqft_than_neighbors -3.716e+04 3131.183 -11.869 0.000 -4.33e+04 -3.1e+04
```

|                       |           |                          |             |
|-----------------------|-----------|--------------------------|-------------|
| <b>Omnibus:</b>       | 21495.431 | <b>Durbin-Watson:</b>    | 1.990       |
| <b>Prob(Omnibus):</b> | 0.000     | <b>Jarque-Bera (JB):</b> | 4815678.728 |
| <b>Skew:</b>          | 4.381     | <b>Prob(JB):</b>         | 0.00        |
| <b>Kurtosis:</b>      | 75.627    | <b>Cond. No.</b>         | 2.47e+08    |

Notes:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [2] The condition number is large, 2.47e+08. This might indicate that there are strong multicollinearity or other numerical problems.

Out[51]: <statsmodels.regression.linear\_model.RegressionResultsWrapper at 0x24d035b2310>



We can now see that our R-squared value for this iteration ended up staying the same at 0.800 and our adjusted R-squared value improved to 0.800 as well after we removed 'sqft\_lot15'. We don't have any remaining issues for the p-values and as we can see our residuals still seem to not be fitting the normality and homoscedasticity assumptions. To address this we can move on to the outlier removal process.

## Outlier Removal

In [52]: #Outlier Removal with the IQR method

```
def find_outliers_IQR(data):
    """Use Tukey's Method of outlier removal AKA InterQuartile-Range Rule
    and return boolean series where True indicates it is an outlier.
    - Calculates the range between the 75% and 25% quartiles
    - Outliers fall outside upper and lower limits, using a threshold of 1.5*IQR the 75

    IQR Range Calculation:
    res = df.describe()
    IQR = res['75%'] - res['25%']
    lower_limit = res['25%'] - 1.5*IQR
    upper_limit = res['75%'] + 1.5*IQR
```

**Args:**  
data (Series, or ndarray): data to test for outliers.

**Returns:**  
[boolean Series]: A True/False for each row use to slice outliers.

**EXAMPLE USE:**

```
>> idx_outs = find_outliers_df(df['AdjustedCompensation'])
>> good_data = df[~idx_outs].copy()
```

function snippet from Flatiron School Phase #2 Py Files.

URL = <https://github.com/flatiron-school/Online-DS-FT-022221-Cohort-Notes/blob/master/02%20-%20Data%20Preparation.ipynb>

'''

```
df_b=data
res= df_b.describe()
```

```
IQR = res['75%'] - res['25%']
lower_limit = res['25%'] - 1.5*IQR
upper_limit = res['75%'] + 1.5*IQR
```

```
idx_outs = (df_b > upper_limit) | (df_b < lower_limit)
```

```
return idx_outs
```

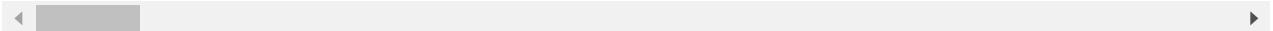
In [53]: *#Making a copy of df\_ohe for the second outlier removal process. Refer to next section.*

```
df_IQR_price = df_ohe.copy()
df_IQR_price
```

Out[53]:

|       | price    | bedrooms | bathrooms | sqft_lot | floors | waterfront | view | condition | grade | sqft_above |
|-------|----------|----------|-----------|----------|--------|------------|------|-----------|-------|------------|
| 0     | 221900.0 | 3        | 1.00      | 5650     | 1.0    | 0.0        | 0.0  | 3         | 7     | 118        |
| 1     | 538000.0 | 3        | 2.25      | 7242     | 2.0    | 0.0        | 0.0  | 3         | 7     | 217        |
| 2     | 180000.0 | 2        | 1.00      | 10000    | 1.0    | 0.0        | 0.0  | 3         | 6     | 77         |
| 3     | 604000.0 | 4        | 3.00      | 5000     | 1.0    | 0.0        | 0.0  | 5         | 7     | 105        |
| 4     | 510000.0 | 3        | 2.00      | 8080     | 1.0    | 0.0        | 0.0  | 3         | 8     | 168        |
| ...   | ...      | ...      | ...       | ...      | ...    | ...        | ...  | ...       | ...   | ...        |
| 21592 | 360000.0 | 3        | 2.50      | 1131     | 3.0    | 0.0        | 0.0  | 3         | 8     | 153        |
| 21593 | 400000.0 | 4        | 2.50      | 5813     | 2.0    | 0.0        | 0.0  | 3         | 8     | 231        |
| 21594 | 402101.0 | 2        | 0.75      | 1350     | 2.0    | 0.0        | 0.0  | 3         | 7     | 102        |
| 21595 | 400000.0 | 3        | 2.50      | 2388     | 2.0    | 0.0        | 0.0  | 3         | 8     | 160        |
| 21596 | 325000.0 | 2        | 0.75      | 1076     | 2.0    | 0.0        | 0.0  | 3         | 7     | 102        |

21597 rows × 85 columns



## Cleaning Outliers From All Numeric Columns

In [54]: `cols_to_check = ['price', 'bedrooms', 'bathrooms', 'sqft_lot', 'grade', 'sqft_above', 'c']`

In [55]: `for col in cols_to_check:`

```
df_ohe = df_ohe[find_outliers_IQR(df_ohe[col]) == False]
df_ohe
```

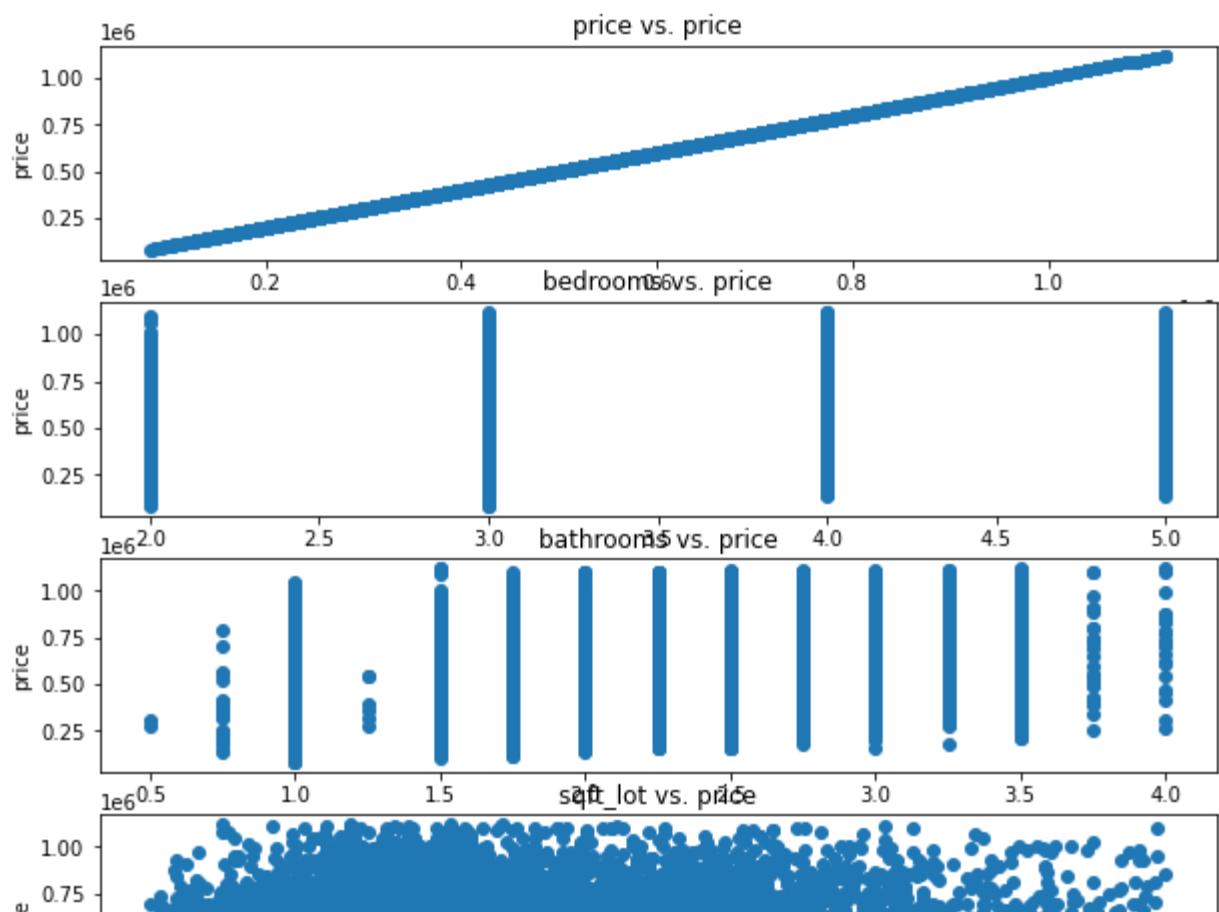
Out[55]:

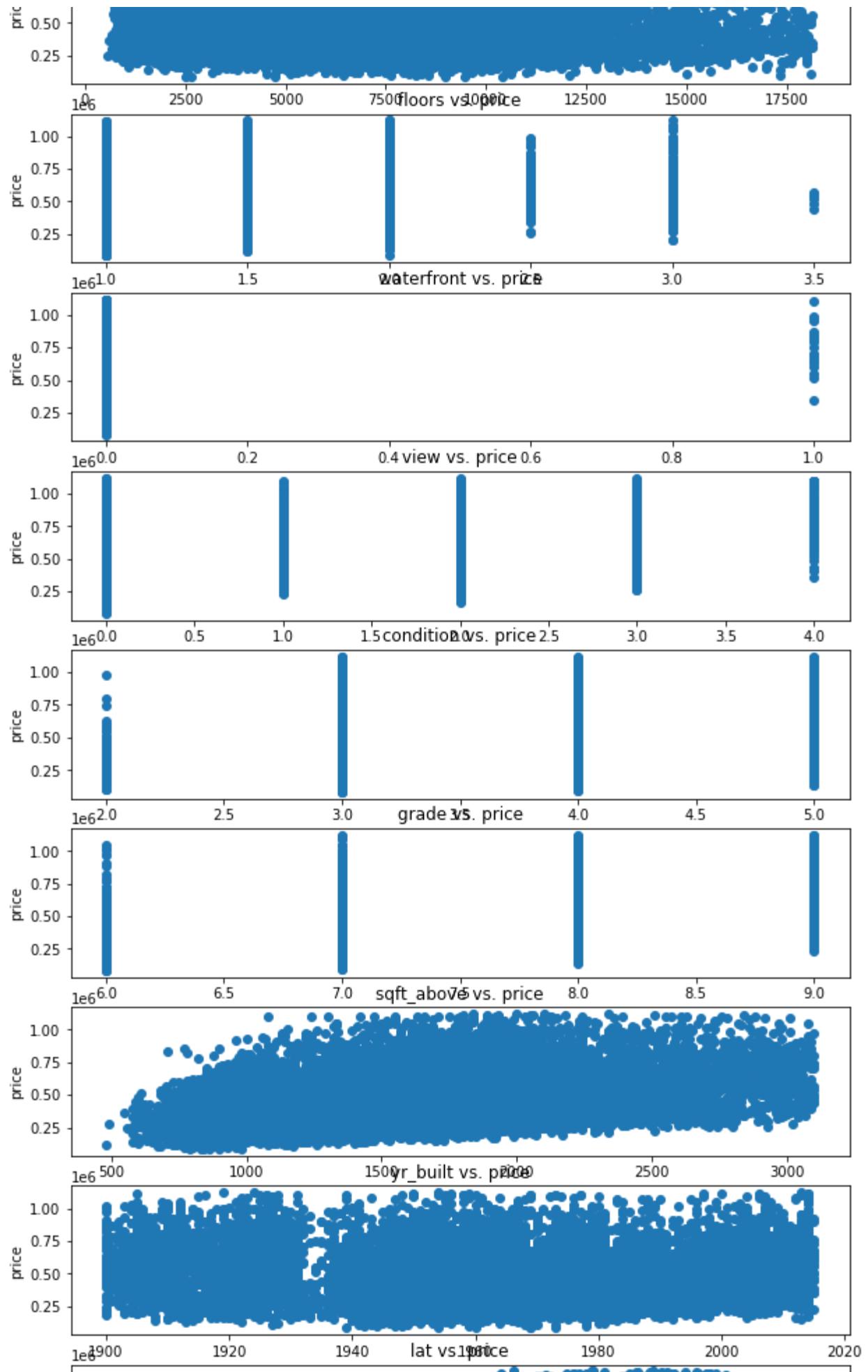
|       | price    | bedrooms | bathrooms | sqft_lot | floors | waterfront | view | condition | grade | sqft_abov |
|-------|----------|----------|-----------|----------|--------|------------|------|-----------|-------|-----------|
| 0     | 221900.0 | 3        | 1.00      | 5650     | 1.0    | 0.0        | 0.0  | 3         | 7     | 118       |
| 1     | 538000.0 | 3        | 2.25      | 7242     | 2.0    | 0.0        | 0.0  | 3         | 7     | 217       |
| 2     | 180000.0 | 2        | 1.00      | 10000    | 1.0    | 0.0        | 0.0  | 3         | 6     | 77        |
| 3     | 604000.0 | 4        | 3.00      | 5000     | 1.0    | 0.0        | 0.0  | 5         | 7     | 105       |
| 4     | 510000.0 | 3        | 2.00      | 8080     | 1.0    | 0.0        | 0.0  | 3         | 8     | 168       |
| ...   | ...      | ...      | ...       | ...      | ...    | ...        | ...  | ...       | ...   | ...       |
| 21592 | 360000.0 | 3        | 2.50      | 1131     | 3.0    | 0.0        | 0.0  | 3         | 8     | 153       |
| 21593 | 400000.0 | 4        | 2.50      | 5813     | 2.0    | 0.0        | 0.0  | 3         | 8     | 231       |
| 21594 | 402101.0 | 2        | 0.75      | 1350     | 2.0    | 0.0        | 0.0  | 3         | 7     | 102       |
| 21595 | 400000.0 | 3        | 2.50      | 2388     | 2.0    | 0.0        | 0.0  | 3         | 8     | 160       |
| 21596 | 325000.0 | 2        | 0.75      | 1076     | 2.0    | 0.0        | 0.0  | 3         | 7     | 102       |

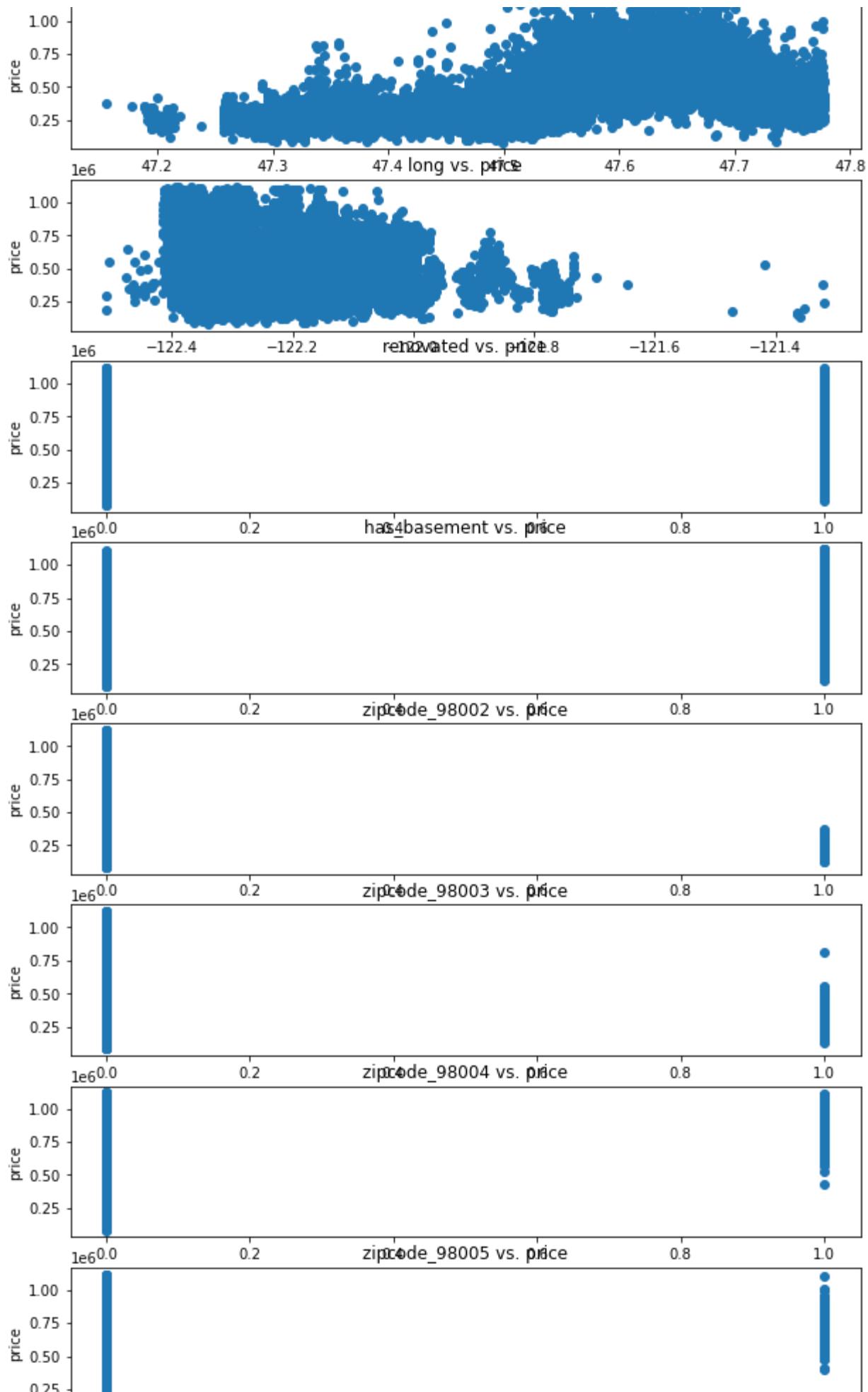
16556 rows × 85 columns

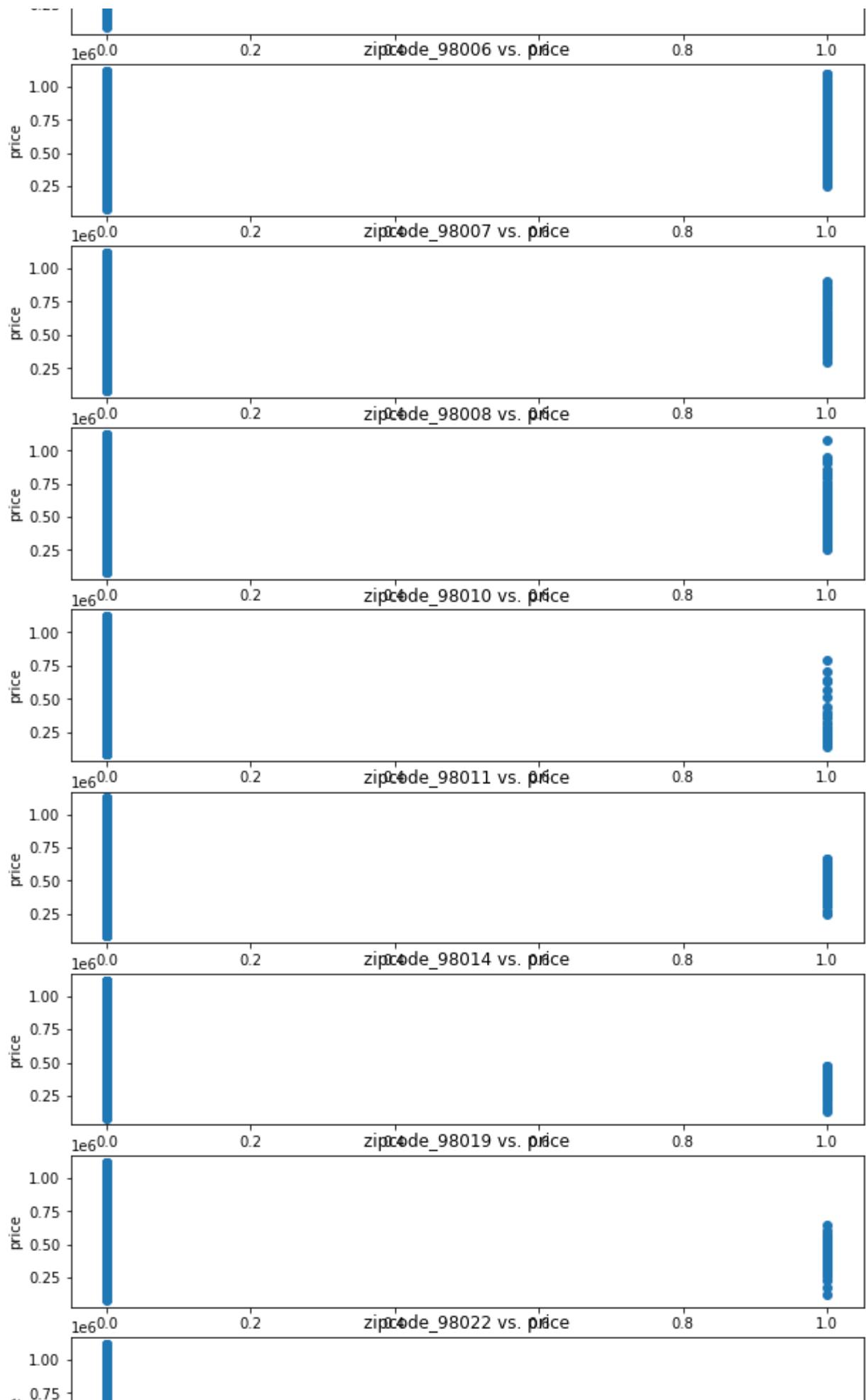
In [56]:

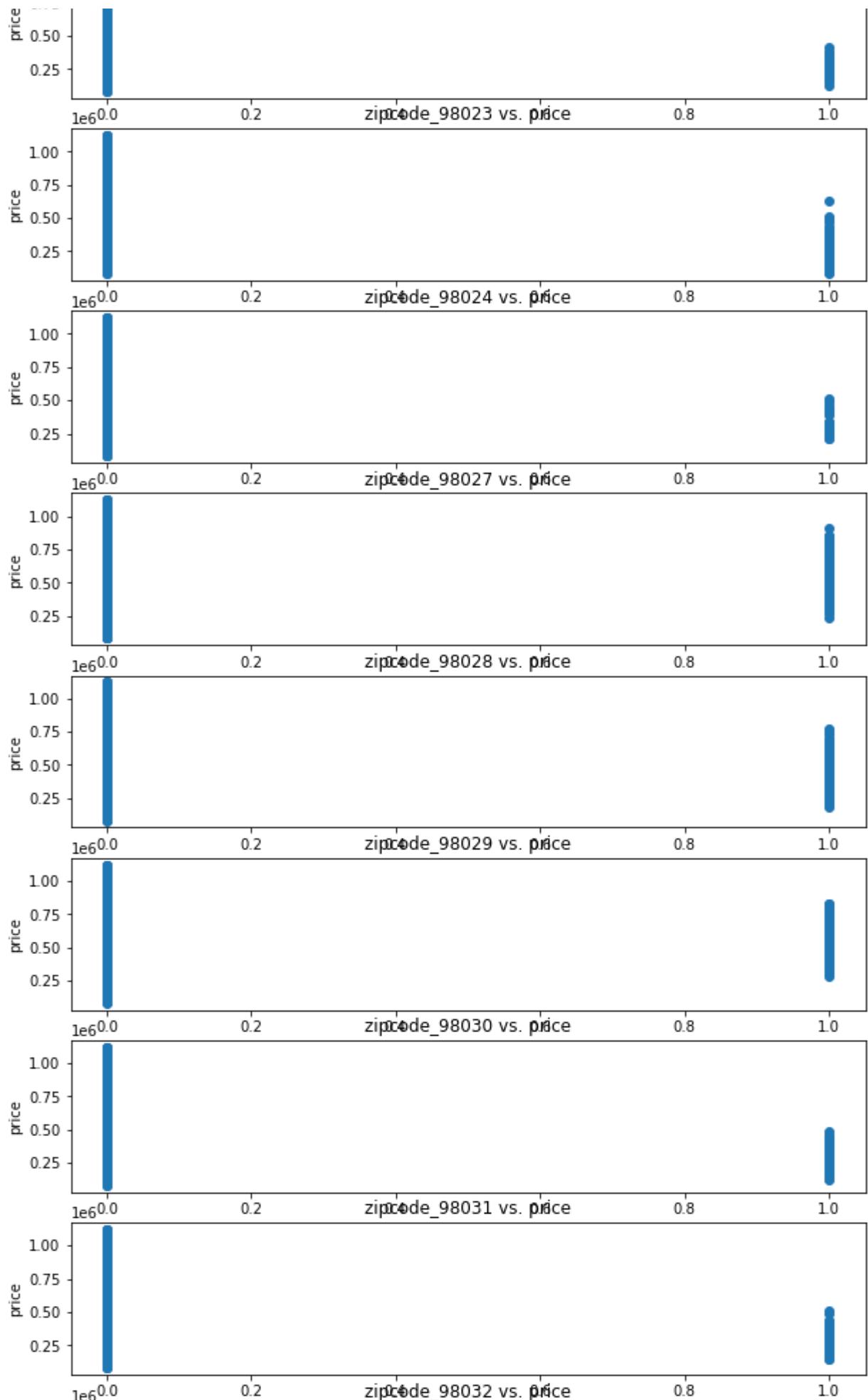
```
#checking the linearity between parameters after outlier removal
plot(df=df_ohe, target='price')
```

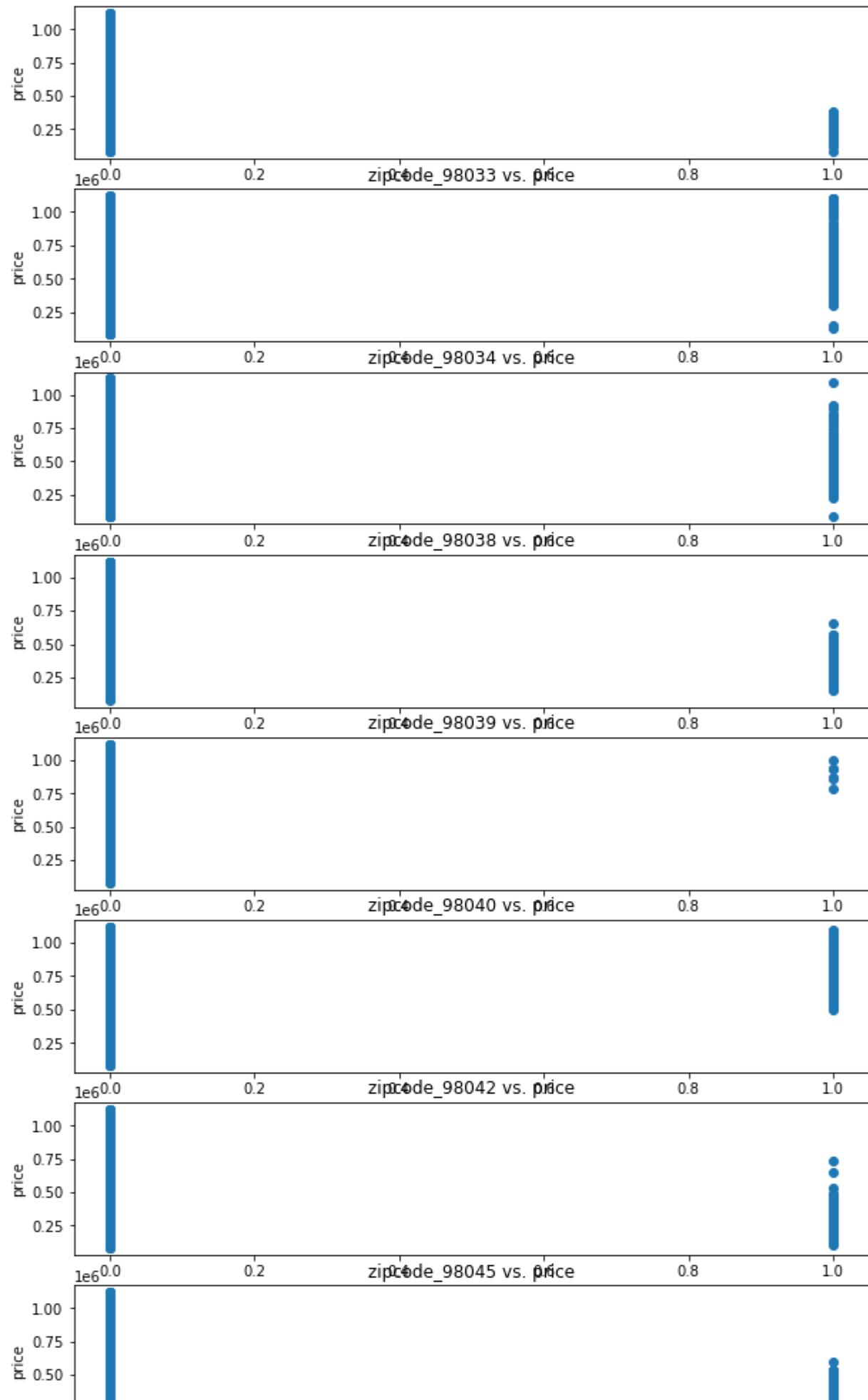


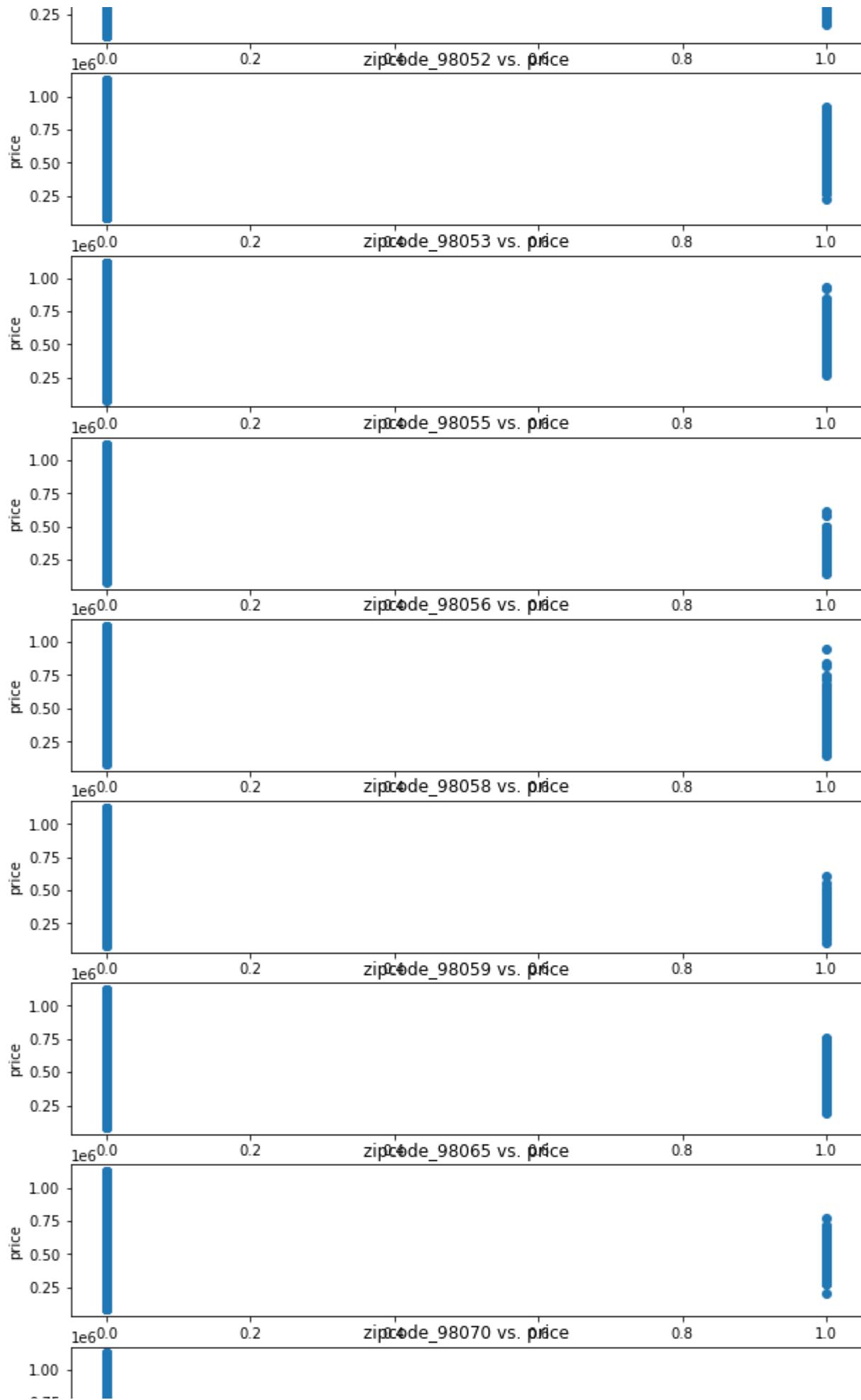


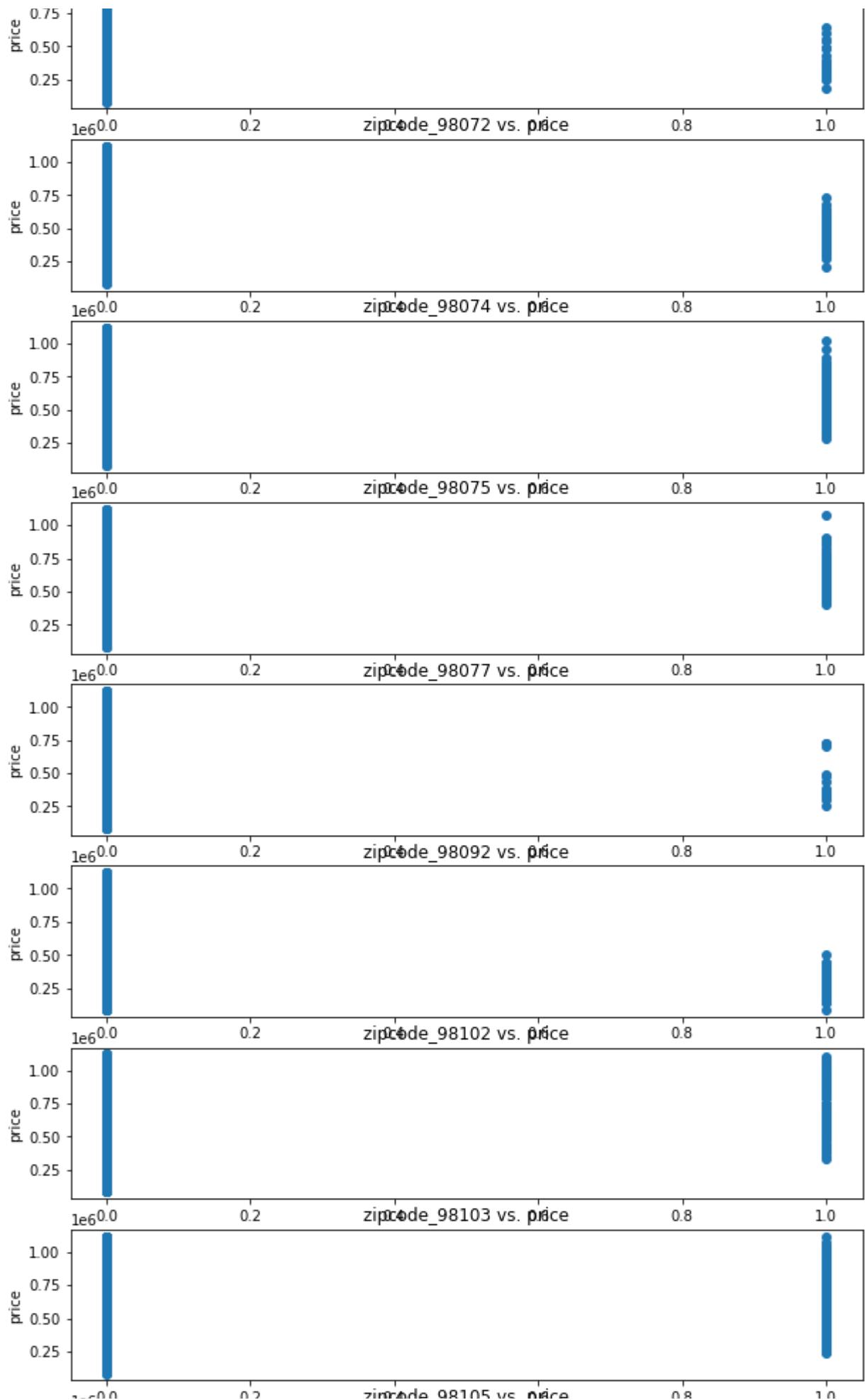


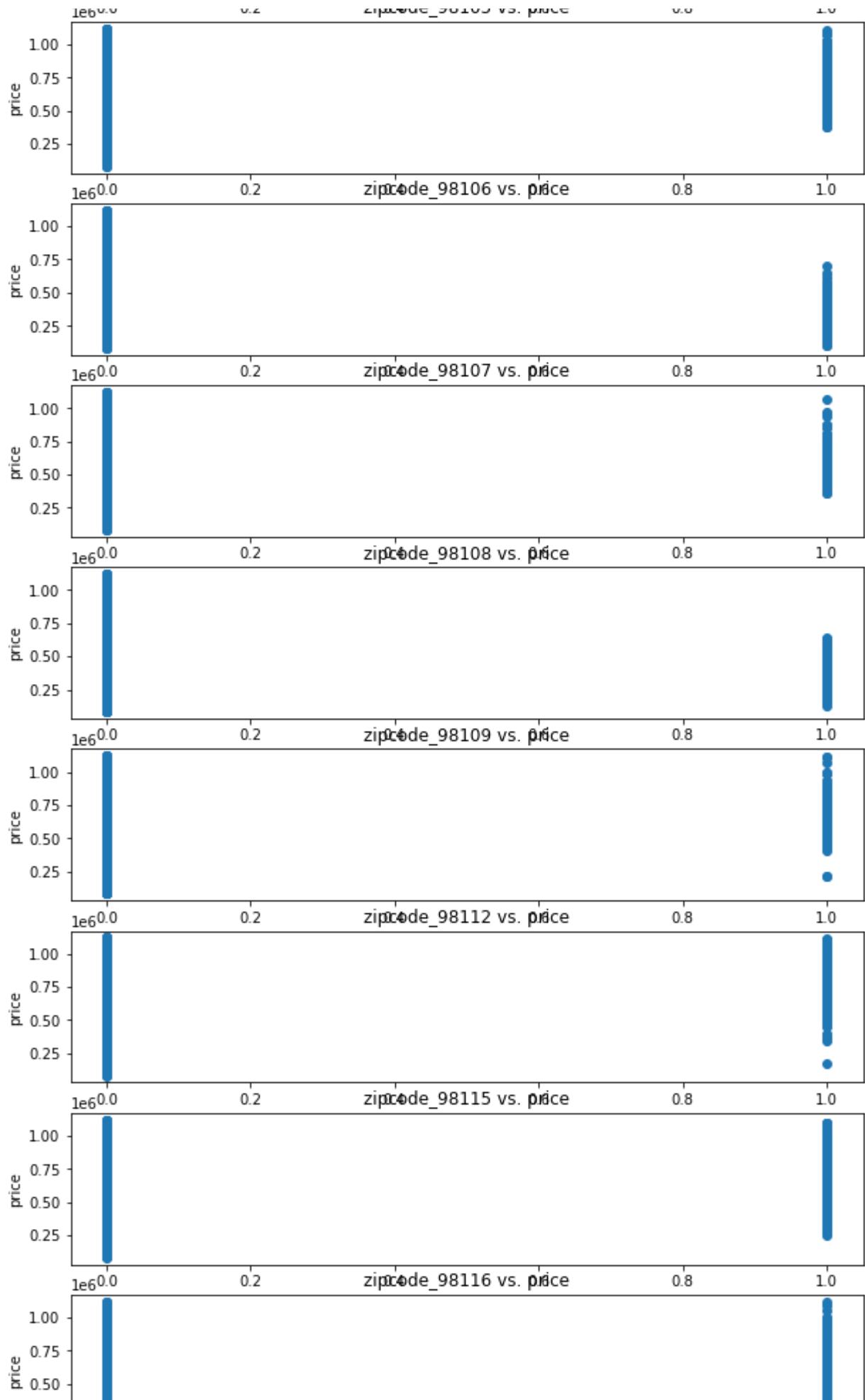


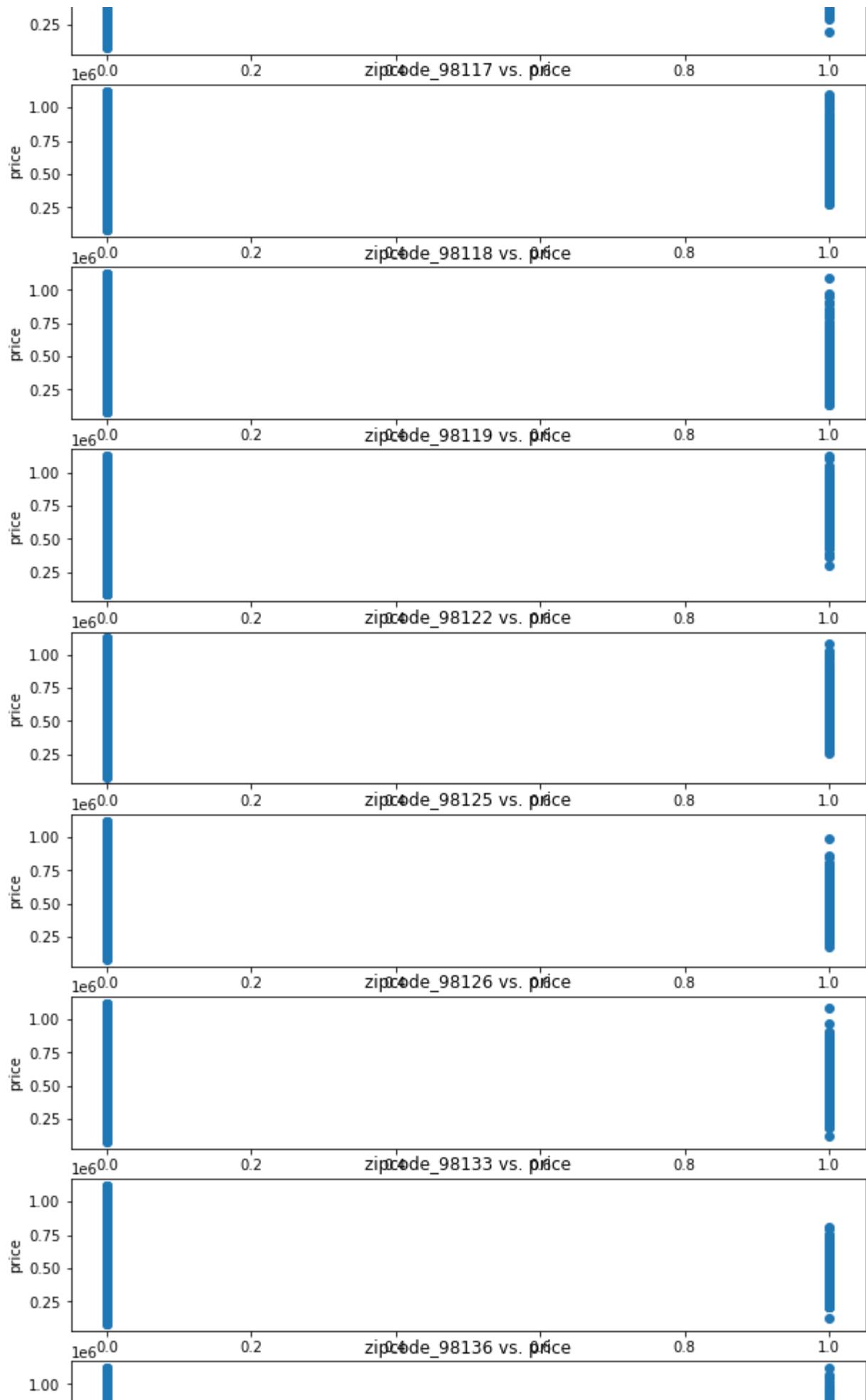


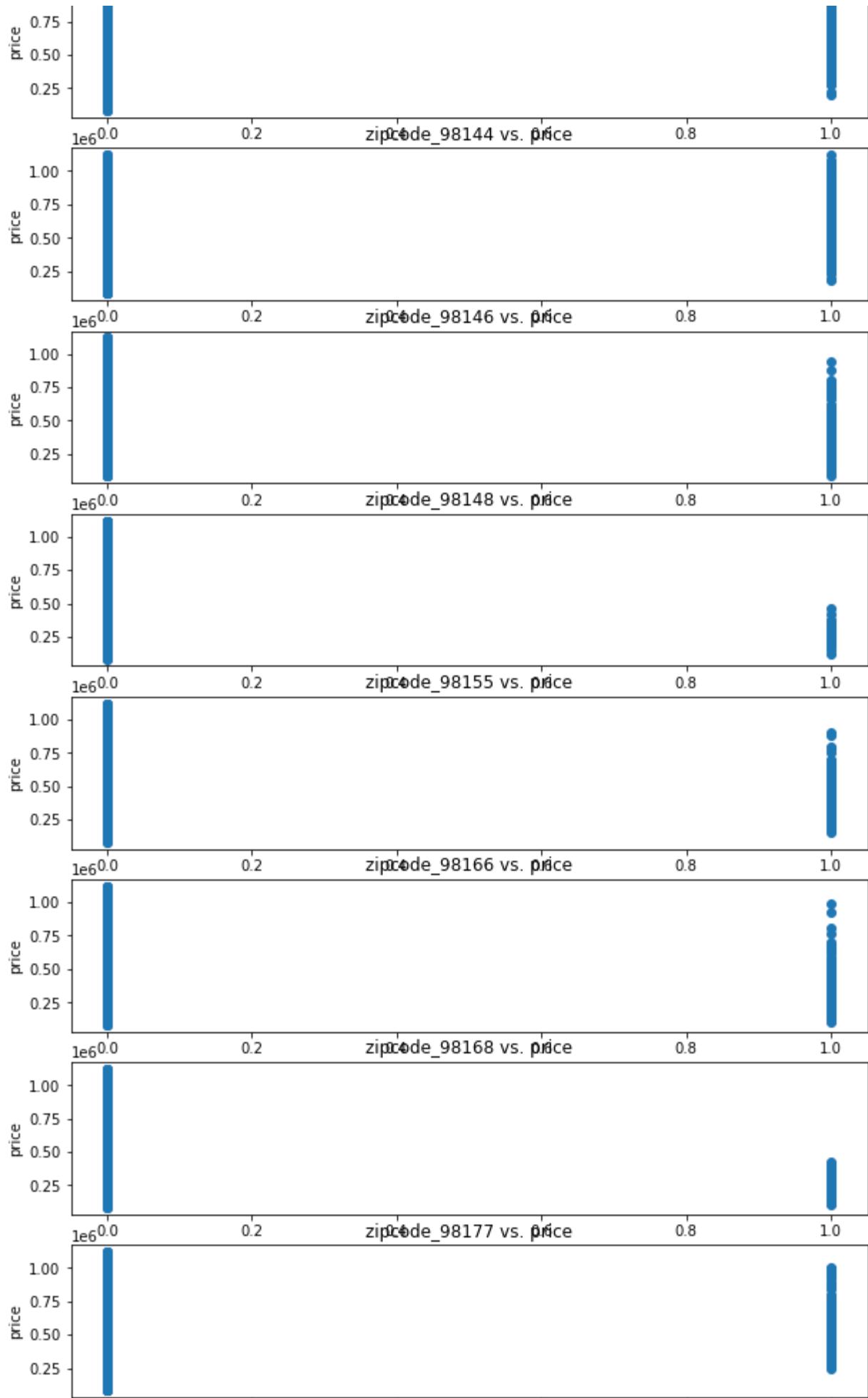


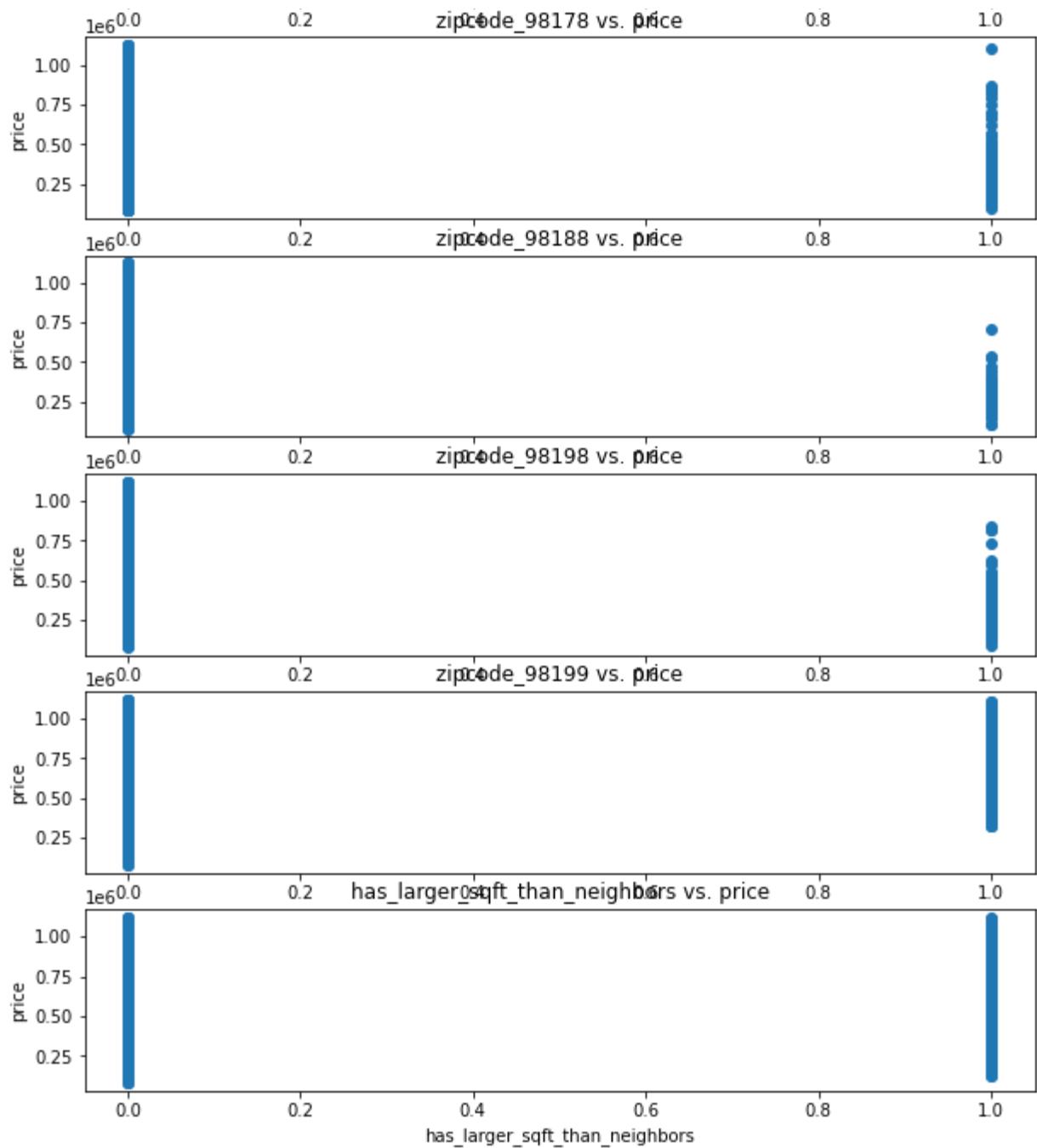












## Model 2.1 - Addressing Outliers

```
In [57]: model = model_lin_reg(df=df_ohe)
```

OLS Regression Results

|                          |                  |                            |             |
|--------------------------|------------------|----------------------------|-------------|
| <b>Dep. Variable:</b>    | price            | <b>R-squared:</b>          | 0.820       |
| <b>Model:</b>            | OLS              | <b>Adj. R-squared:</b>     | 0.819       |
| <b>Method:</b>           | Least Squares    | <b>F-statistic:</b>        | 890.9       |
| <b>Date:</b>             | Sat, 01 May 2021 | <b>Prob (F-statistic):</b> | 0.00        |
| <b>Time:</b>             | 16:50:36         | <b>Log-Likelihood:</b>     | -2.1045e+05 |
| <b>No. Observations:</b> | 16556            | <b>AIC:</b>                | 4.211e+05   |
| <b>Df Residuals:</b>     | 16471            | <b>BIC:</b>                | 4.217e+05   |

**Df Model:** 84**Covariance Type:** nonrobust

|  |                      | <b>coef</b> | <b>std err</b> | <b>t</b> | <b>P&gt; t </b> | <b>[0.025</b> | <b>0.975]</b> |
|--|----------------------|-------------|----------------|----------|-----------------|---------------|---------------|
|  | <b>Intercept</b>     | -1.459e+06  | 4.28e+06       | -0.341   | 0.733           | -9.85e+06     | 6.94e+06      |
|  | <b>bedrooms</b>      | 2970.0115   | 1032.190       | 2.877    | 0.004           | 946.808       | 4993.215      |
|  | <b>bathrooms</b>     | 2.089e+04   | 1587.628       | 13.158   | 0.000           | 1.78e+04      | 2.4e+04       |
|  | <b>sqft_lot</b>      | 1.9798      | 0.268          | 7.380    | 0.000           | 1.454         | 2.506         |
|  | <b>floors</b>        | -2.331e+04  | 1931.424       | -12.066  | 0.000           | -2.71e+04     | -1.95e+04     |
|  | <b>waterfront</b>    | 2.397e+05   | 1.88e+04       | 12.747   | 0.000           | 2.03e+05      | 2.77e+05      |
|  | <b>view</b>          | 4.034e+04   | 1180.302       | 34.180   | 0.000           | 3.8e+04       | 4.27e+04      |
|  | <b>condition</b>     | 2.517e+04   | 1104.208       | 22.795   | 0.000           | 2.3e+04       | 2.73e+04      |
|  | <b>grade</b>         | 4.6e+04     | 1243.293       | 37.000   | 0.000           | 4.36e+04      | 4.84e+04      |
|  | <b>sqft_above</b>    | 139.7277    | 2.321          | 60.207   | 0.000           | 135.179       | 144.277       |
|  | <b>yr_built</b>      | -591.2032   | 38.117         | -15.510  | 0.000           | -665.916      | -516.490      |
|  | <b>lat</b>           | -4840.7162  | 4.02e+04       | -0.120   | 0.904           | -8.36e+04     | 7.39e+04      |
|  | <b>long</b>          | -1.989e+04  | 3.25e+04       | -0.613   | 0.540           | -8.35e+04     | 4.37e+04      |
|  | <b>renovated</b>     | 3.521e+04   | 3896.615       | 9.037    | 0.000           | 2.76e+04      | 4.28e+04      |
|  | <b>has_basement</b>  | 4.734e+04   | 1730.499       | 27.359   | 0.000           | 4.4e+04       | 5.07e+04      |
|  | <b>zipcode_98002</b> | 6751.8751   | 8013.648       | 0.843    | 0.399           | -8955.740     | 2.25e+04      |
|  | <b>zipcode_98003</b> | -2498.5176  | 7128.191       | -0.351   | 0.726           | -1.65e+04     | 1.15e+04      |
|  | <b>zipcode_98004</b> | 5.278e+05   | 1.51e+04       | 35.017   | 0.000           | 4.98e+05      | 5.57e+05      |
|  | <b>zipcode_98005</b> | 3.347e+05   | 1.54e+04       | 21.725   | 0.000           | 3.05e+05      | 3.65e+05      |
|  | <b>zipcode_98006</b> | 2.806e+05   | 1.28e+04       | 21.865   | 0.000           | 2.55e+05      | 3.06e+05      |
|  | <b>zipcode_98007</b> | 2.48e+05    | 1.58e+04       | 15.711   | 0.000           | 2.17e+05      | 2.79e+05      |
|  | <b>zipcode_98008</b> | 2.408e+05   | 1.52e+04       | 15.812   | 0.000           | 2.11e+05      | 2.71e+05      |
|  | <b>zipcode_98010</b> | 8.48e+04    | 1.61e+04       | 5.271    | 0.000           | 5.33e+04      | 1.16e+05      |
|  | <b>zipcode_98011</b> | 1.502e+05   | 1.98e+04       | 7.596    | 0.000           | 1.11e+05      | 1.89e+05      |
|  | <b>zipcode_98014</b> | 1.139e+05   | 2.52e+04       | 4.518    | 0.000           | 6.45e+04      | 1.63e+05      |
|  | <b>zipcode_98019</b> | 1.076e+05   | 2.2e+04        | 4.895    | 0.000           | 6.45e+04      | 1.51e+05      |
|  | <b>zipcode_98022</b> | 1843.3567   | 1.28e+04       | 0.144    | 0.886           | -2.33e+04     | 2.69e+04      |
|  | <b>zipcode_98023</b> | -1.568e+04  | 6912.313       | -2.269   | 0.023           | -2.92e+04     | -2136.085     |
|  | <b>zipcode_98024</b> | 1.319e+05   | 2.4e+04        | 5.484    | 0.000           | 8.47e+04      | 1.79e+05      |
|  | <b>zipcode_98027</b> | 2.319e+05   | 1.41e+04       | 16.451   | 0.000           | 2.04e+05      | 2.59e+05      |
|  | <b>zipcode_98028</b> | 1.39e+05    | 1.92e+04       | 7.226    | 0.000           | 1.01e+05      | 1.77e+05      |

|                      | final_notebook-Copy1 |          |        |       |           |           |  |
|----------------------|----------------------|----------|--------|-------|-----------|-----------|--|
| <b>zipcode_98029</b> | 2.395e+05            | 1.54e+04 | 15.582 | 0.000 | 2.09e+05  | 2.7e+05   |  |
| <b>zipcode_98030</b> | 7538.6621            | 8010.111 | 0.941  | 0.347 | -8162.021 | 2.32e+04  |  |
| <b>zipcode_98031</b> | 1.367e+04            | 8489.042 | 1.610  | 0.107 | -2970.898 | 3.03e+04  |  |
| <b>zipcode_98032</b> | -7011.8975           | 9410.308 | -0.745 | 0.456 | -2.55e+04 | 1.14e+04  |  |
| <b>zipcode_98033</b> | 3.175e+05            | 1.67e+04 | 19.007 | 0.000 | 2.85e+05  | 3.5e+05   |  |
| <b>zipcode_98034</b> | 1.866e+05            | 1.78e+04 | 10.467 | 0.000 | 1.52e+05  | 2.22e+05  |  |
| <b>zipcode_98038</b> | 4.712e+04            | 1.01e+04 | 4.652  | 0.000 | 2.73e+04  | 6.7e+04   |  |
| <b>zipcode_98039</b> | 6.561e+05            | 3.57e+04 | 18.399 | 0.000 | 5.86e+05  | 7.26e+05  |  |
| <b>zipcode_98040</b> | 4.384e+05            | 1.34e+04 | 32.602 | 0.000 | 4.12e+05  | 4.65e+05  |  |
| <b>zipcode_98042</b> | 1.652e+04            | 8413.643 | 1.963  | 0.050 | 27.240    | 3.3e+04   |  |
| <b>zipcode_98045</b> | 1.03e+05             | 1.99e+04 | 5.185  | 0.000 | 6.41e+04  | 1.42e+05  |  |
| <b>zipcode_98052</b> | 2.576e+05            | 1.7e+04  | 15.184 | 0.000 | 2.24e+05  | 2.91e+05  |  |
| <b>zipcode_98053</b> | 2.677e+05            | 1.94e+04 | 13.793 | 0.000 | 2.3e+05   | 3.06e+05  |  |
| <b>zipcode_98055</b> | 4.165e+04            | 9647.913 | 4.317  | 0.000 | 2.27e+04  | 6.06e+04  |  |
| <b>zipcode_98056</b> | 1.02e+05             | 1.09e+04 | 9.359  | 0.000 | 8.06e+04  | 1.23e+05  |  |
| <b>zipcode_98058</b> | 3.764e+04            | 9421.517 | 3.995  | 0.000 | 1.92e+04  | 5.61e+04  |  |
| <b>zipcode_98059</b> | 9.328e+04            | 1.08e+04 | 8.662  | 0.000 | 7.22e+04  | 1.14e+05  |  |
| <b>zipcode_98065</b> | 1.486e+05            | 1.79e+04 | 8.310  | 0.000 | 1.14e+05  | 1.84e+05  |  |
| <b>zipcode_98070</b> | 4.969e+04            | 1.84e+04 | 2.703  | 0.007 | 1.37e+04  | 8.57e+04  |  |
| <b>zipcode_98072</b> | 1.613e+05            | 2.03e+04 | 7.967  | 0.000 | 1.22e+05  | 2.01e+05  |  |
| <b>zipcode_98074</b> | 2.159e+05            | 1.66e+04 | 13.019 | 0.000 | 1.83e+05  | 2.48e+05  |  |
| <b>zipcode_98075</b> | 2.516e+05            | 1.69e+04 | 14.896 | 0.000 | 2.19e+05  | 2.85e+05  |  |
| <b>zipcode_98077</b> | 1.486e+05            | 2.49e+04 | 5.980  | 0.000 | 9.99e+04  | 1.97e+05  |  |
| <b>zipcode_98092</b> | -1.782e+04           | 7594.534 | -2.346 | 0.019 | -3.27e+04 | -2929.013 |  |
| <b>zipcode_98102</b> | 4.117e+05            | 1.67e+04 | 24.706 | 0.000 | 3.79e+05  | 4.44e+05  |  |
| <b>zipcode_98103</b> | 3.265e+05            | 1.6e+04  | 20.422 | 0.000 | 2.95e+05  | 3.58e+05  |  |
| <b>zipcode_98105</b> | 3.773e+05            | 1.64e+04 | 22.969 | 0.000 | 3.45e+05  | 4.09e+05  |  |
| <b>zipcode_98106</b> | 1.215e+05            | 1.15e+04 | 10.595 | 0.000 | 9.91e+04  | 1.44e+05  |  |
| <b>zipcode_98107</b> | 3.233e+05            | 1.63e+04 | 19.820 | 0.000 | 2.91e+05  | 3.55e+05  |  |
| <b>zipcode_98108</b> | 1.229e+05            | 1.24e+04 | 9.887  | 0.000 | 9.85e+04  | 1.47e+05  |  |
| <b>zipcode_98109</b> | 4.129e+05            | 1.68e+04 | 24.634 | 0.000 | 3.8e+05   | 4.46e+05  |  |
| <b>zipcode_98112</b> | 4.353e+05            | 1.52e+04 | 28.561 | 0.000 | 4.05e+05  | 4.65e+05  |  |
| <b>zipcode_98115</b> | 3.232e+05            | 1.62e+04 | 19.899 | 0.000 | 2.91e+05  | 3.55e+05  |  |
| <b>zipcode_98116</b> | 2.961e+05            | 1.3e+04  | 22.728 | 0.000 | 2.71e+05  | 3.22e+05  |  |
| <b>zipcode_98117</b> | 3.126e+05            | 1.65e+04 | 18.945 | 0.000 | 2.8e+05   | 3.45e+05  |  |

|                                       |            |          |        |       |           |           |
|---------------------------------------|------------|----------|--------|-------|-----------|-----------|
| <b>zipcode_98118</b>                  | 1.693e+05  | 1.13e+04 | 15.011 | 0.000 | 1.47e+05  | 1.91e+05  |
| <b>zipcode_98119</b>                  | 4.056e+05  | 1.59e+04 | 25.484 | 0.000 | 3.74e+05  | 4.37e+05  |
| <b>zipcode_98122</b>                  | 3.107e+05  | 1.41e+04 | 22.080 | 0.000 | 2.83e+05  | 3.38e+05  |
| <b>zipcode_98125</b>                  | 1.964e+05  | 1.75e+04 | 11.220 | 0.000 | 1.62e+05  | 2.31e+05  |
| <b>zipcode_98126</b>                  | 1.911e+05  | 1.18e+04 | 16.139 | 0.000 | 1.68e+05  | 2.14e+05  |
| <b>zipcode_98133</b>                  | 1.55e+05   | 1.81e+04 | 8.553  | 0.000 | 1.2e+05   | 1.91e+05  |
| <b>zipcode_98136</b>                  | 2.547e+05  | 1.21e+04 | 21.058 | 0.000 | 2.31e+05  | 2.78e+05  |
| <b>zipcode_98144</b>                  | 2.453e+05  | 1.31e+04 | 18.764 | 0.000 | 2.2e+05   | 2.71e+05  |
| <b>zipcode_98146</b>                  | 1.054e+05  | 1.07e+04 | 9.876  | 0.000 | 8.45e+04  | 1.26e+05  |
| <b>zipcode_98148</b>                  | 4.602e+04  | 1.33e+04 | 3.473  | 0.001 | 2e+04     | 7.2e+04   |
| <b>zipcode_98155</b>                  | 1.424e+05  | 1.88e+04 | 7.555  | 0.000 | 1.05e+05  | 1.79e+05  |
| <b>zipcode_98166</b>                  | 8.588e+04  | 9838.761 | 8.729  | 0.000 | 6.66e+04  | 1.05e+05  |
| <b>zipcode_98168</b>                  | 5.147e+04  | 1.03e+04 | 4.997  | 0.000 | 3.13e+04  | 7.17e+04  |
| <b>zipcode_98177</b>                  | 2.054e+05  | 1.89e+04 | 10.855 | 0.000 | 1.68e+05  | 2.43e+05  |
| <b>zipcode_98178</b>                  | 5.611e+04  | 1.05e+04 | 5.353  | 0.000 | 3.56e+04  | 7.67e+04  |
| <b>zipcode_98188</b>                  | 3.822e+04  | 1.04e+04 | 3.669  | 0.000 | 1.78e+04  | 5.86e+04  |
| <b>zipcode_98198</b>                  | 2.03e+04   | 7895.317 | 2.571  | 0.010 | 4825.827  | 3.58e+04  |
| <b>zipcode_98199</b>                  | 3.629e+05  | 1.57e+04 | 23.181 | 0.000 | 3.32e+05  | 3.94e+05  |
| <b>has_larger_sqft_than_neighbors</b> | -1.729e+04 | 1801.693 | -9.599 | 0.000 | -2.08e+04 | -1.38e+04 |

**Omnibus:** 1493.906    **Durbin-Watson:** 1.993

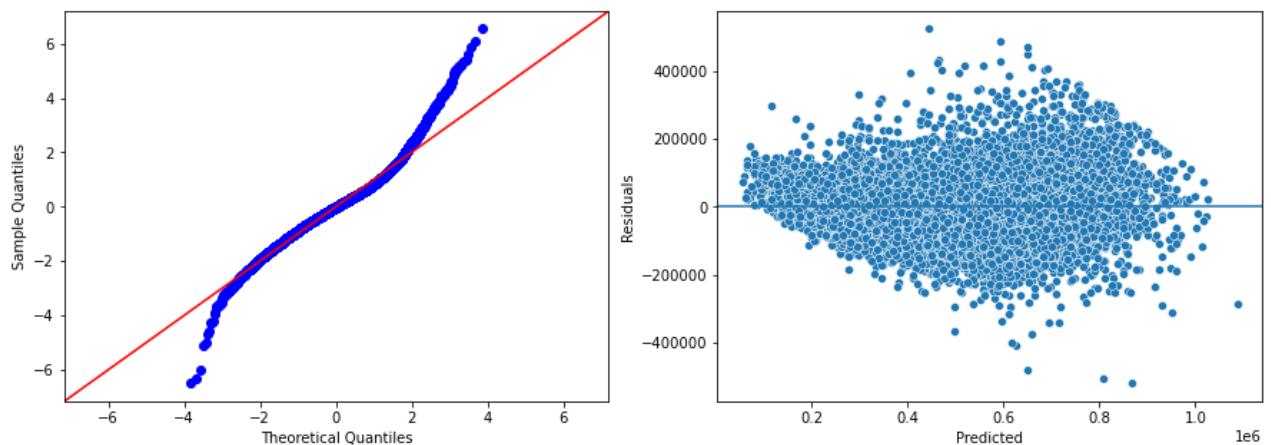
**Prob(Omnibus):** 0.000    **Jarque-Bera (JB):** 5130.230

**Skew:** 0.437    **Prob(JB):** 0.00

**Kurtosis:** 5.583    **Cond. No.** 5.61e+07

#### Notes:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [2] The condition number is large, 5.61e+07. This might indicate that there are strong multicollinearity or other numerical problems.



In [58]: `model.pvalues[model.pvalues>0.05]`

Out[58]:

| Intercept     | 0.733445 |
|---------------|----------|
| lat           | 0.904094 |
| long          | 0.540053 |
| zipcode_98002 | 0.399494 |
| zipcode_98003 | 0.725959 |
| zipcode_98022 | 0.885568 |
| zipcode_98030 | 0.346645 |
| zipcode_98031 | 0.107387 |
| zipcode_98032 | 0.456204 |
| dtype:        | float64  |

Our model's R-squared and adjusted R-squared values improved to 0.820 and 0.819 respectively compared to our baseline model's 0.800 and 0.799. From our p-values, we can see that the latitude and longitude values are insignificant which means that we can go ahead and drop these coefficients. Our residuals are also looking more normal and homoscedastic compared to the baseline model's residuals.

In [59]: `df_ohe.drop(['lat', 'long'], axis=1, inplace=True)`

In [60]: `model = model_lin_reg(df=df_ohe)`

### OLS Regression Results

| Dep. Variable:    | price            | R-squared:          | 0.820       |
|-------------------|------------------|---------------------|-------------|
| Model:            | OLS              | Adj. R-squared:     | 0.819       |
| Method:           | Least Squares    | F-statistic:        | 912.7       |
| Date:             | Sat, 01 May 2021 | Prob (F-statistic): | 0.00        |
| Time:             | 16:50:37         | Log-Likelihood:     | -2.1045e+05 |
| No. Observations: | 16556            | AIC:                | 4.211e+05   |
| Df Residuals:     | 16473            | BIC:                | 4.217e+05   |
| Df Model:         | 82               |                     |             |
| Covariance Type:  | nonrobust        |                     |             |

|           | coef      | std err  | t      | P> t  | [0.025 | 0.975]   |
|-----------|-----------|----------|--------|-------|--------|----------|
| Intercept | 7.453e+05 | 7.43e+04 | 10.028 | 0.000 | 6e+05  | 8.91e+05 |

|                      |            |          |         |       |           |           |
|----------------------|------------|----------|---------|-------|-----------|-----------|
| <b>bedrooms</b>      | 2967.3423  | 1032.116 | 2.875   | 0.004 | 944.283   | 4990.401  |
| <b>bathrooms</b>     | 2.089e+04  | 1587.508 | 13.159  | 0.000 | 1.78e+04  | 2.4e+04   |
| <b>sqft_lot</b>      | 1.9769     | 0.268    | 7.371   | 0.000 | 1.451     | 2.503     |
| <b>floors</b>        | -2.329e+04 | 1930.880 | -12.061 | 0.000 | -2.71e+04 | -1.95e+04 |
| <b>waterfront</b>    | 2.398e+05  | 1.88e+04 | 12.752  | 0.000 | 2.03e+05  | 2.77e+05  |
| <b>view</b>          | 4.036e+04  | 1179.978 | 34.201  | 0.000 | 3.8e+04   | 4.27e+04  |
| <b>condition</b>     | 2.516e+04  | 1103.866 | 22.796  | 0.000 | 2.3e+04   | 2.73e+04  |
| <b>grade</b>         | 4.602e+04  | 1242.410 | 37.043  | 0.000 | 4.36e+04  | 4.85e+04  |
| <b>sqft_above</b>    | 139.7018   | 2.320    | 60.210  | 0.000 | 135.154   | 144.250   |
| <b>yr_built</b>      | -591.9980  | 38.092   | -15.541 | 0.000 | -666.663  | -517.333  |
| <b>renovated</b>     | 3.518e+04  | 3896.042 | 9.030   | 0.000 | 2.75e+04  | 4.28e+04  |
| <b>has_basement</b>  | 4.735e+04  | 1730.358 | 27.365  | 0.000 | 4.4e+04   | 5.07e+04  |
| <b>zipcode_98002</b> | 5621.7555  | 7796.170 | 0.721   | 0.471 | -9659.580 | 2.09e+04  |
| <b>zipcode_98003</b> | -1734.1376 | 7011.609 | -0.247  | 0.805 | -1.55e+04 | 1.2e+04   |
| <b>zipcode_98004</b> | 5.249e+05  | 8553.918 | 61.369  | 0.000 | 5.08e+05  | 5.42e+05  |
| <b>zipcode_98005</b> | 3.313e+05  | 9253.837 | 35.803  | 0.000 | 3.13e+05  | 3.49e+05  |
| <b>zipcode_98006</b> | 2.77e+05   | 6768.602 | 40.926  | 0.000 | 2.64e+05  | 2.9e+05   |
| <b>zipcode_98007</b> | 2.439e+05  | 9008.742 | 27.077  | 0.000 | 2.26e+05  | 2.62e+05  |
| <b>zipcode_98008</b> | 2.362e+05  | 7104.393 | 33.253  | 0.000 | 2.22e+05  | 2.5e+05   |
| <b>zipcode_98010</b> | 7.967e+04  | 1.38e+04 | 5.785   | 0.000 | 5.27e+04  | 1.07e+05  |
| <b>zipcode_98011</b> | 1.466e+05  | 7965.911 | 18.408  | 0.000 | 1.31e+05  | 1.62e+05  |
| <b>zipcode_98014</b> | 1.033e+05  | 1.35e+04 | 7.655   | 0.000 | 7.69e+04  | 1.3e+05   |
| <b>zipcode_98019</b> | 9.954e+04  | 8480.295 | 11.737  | 0.000 | 8.29e+04  | 1.16e+05  |
| <b>zipcode_98022</b> | -3106.2087 | 8379.925 | -0.371  | 0.711 | -1.95e+04 | 1.33e+04  |
| <b>zipcode_98023</b> | -1.38e+04  | 6195.130 | -2.227  | 0.026 | -2.59e+04 | -1654.491 |
| <b>zipcode_98024</b> | 1.233e+05  | 1.78e+04 | 6.908   | 0.000 | 8.83e+04  | 1.58e+05  |
| <b>zipcode_98027</b> | 2.265e+05  | 7533.833 | 30.059  | 0.000 | 2.12e+05  | 2.41e+05  |
| <b>zipcode_98028</b> | 1.364e+05  | 7089.711 | 19.235  | 0.000 | 1.22e+05  | 1.5e+05   |
| <b>zipcode_98029</b> | 2.33e+05   | 7011.710 | 33.233  | 0.000 | 2.19e+05  | 2.47e+05  |
| <b>zipcode_98030</b> | 5647.8650  | 7142.783 | 0.791   | 0.429 | -8352.761 | 1.96e+04  |
| <b>zipcode_98031</b> | 1.16e+04   | 7060.845 | 1.643   | 0.100 | -2241.657 | 2.54e+04  |
| <b>zipcode_98032</b> | -7118.2546 | 9084.315 | -0.784  | 0.433 | -2.49e+04 | 1.07e+04  |
| <b>zipcode_98033</b> | 3.14e+05   | 6664.977 | 47.108  | 0.000 | 3.01e+05  | 3.27e+05  |
| <b>zipcode_98034</b> | 1.833e+05  | 6081.409 | 30.142  | 0.000 | 1.71e+05  | 1.95e+05  |
| <b>zipcode_98038</b> | 4.218e+04  | 6123.375 | 6.889   | 0.000 | 3.02e+04  | 5.42e+04  |

|                      |            |          |        |       |           |           |
|----------------------|------------|----------|--------|-------|-----------|-----------|
| <b>zipcode_98039</b> | 6.537e+05  | 3.32e+04 | 19.684 | 0.000 | 5.89e+05  | 7.19e+05  |
| <b>zipcode_98040</b> | 4.362e+05  | 8660.795 | 50.370 | 0.000 | 4.19e+05  | 4.53e+05  |
| <b>zipcode_98042</b> | 1.316e+04  | 6237.511 | 2.111  | 0.035 | 938.149   | 2.54e+04  |
| <b>zipcode_98045</b> | 9.222e+04  | 8411.472 | 10.963 | 0.000 | 7.57e+04  | 1.09e+05  |
| <b>zipcode_98052</b> | 2.528e+05  | 6211.997 | 40.696 | 0.000 | 2.41e+05  | 2.65e+05  |
| <b>zipcode_98053</b> | 2.61e+05   | 7611.304 | 34.286 | 0.000 | 2.46e+05  | 2.76e+05  |
| <b>zipcode_98055</b> | 3.957e+04  | 7202.509 | 5.494  | 0.000 | 2.55e+04  | 5.37e+04  |
| <b>zipcode_98056</b> | 9.924e+04  | 6610.285 | 15.013 | 0.000 | 8.63e+04  | 1.12e+05  |
| <b>zipcode_98058</b> | 3.462e+04  | 6422.863 | 5.391  | 0.000 | 2.2e+04   | 4.72e+04  |
| <b>zipcode_98059</b> | 8.987e+04  | 6561.505 | 13.696 | 0.000 | 7.7e+04   | 1.03e+05  |
| <b>zipcode_98065</b> | 1.395e+05  | 7348.861 | 18.986 | 0.000 | 1.25e+05  | 1.54e+05  |
| <b>zipcode_98070</b> | 5.291e+04  | 1.7e+04  | 3.114  | 0.002 | 1.96e+04  | 8.62e+04  |
| <b>zipcode_98072</b> | 1.567e+05  | 8384.514 | 18.688 | 0.000 | 1.4e+05   | 1.73e+05  |
| <b>zipcode_98074</b> | 2.098e+05  | 6994.896 | 29.998 | 0.000 | 1.96e+05  | 2.24e+05  |
| <b>zipcode_98075</b> | 2.455e+05  | 9339.484 | 26.284 | 0.000 | 2.27e+05  | 2.64e+05  |
| <b>zipcode_98077</b> | 1.426e+05  | 1.6e+04  | 8.928  | 0.000 | 1.11e+05  | 1.74e+05  |
| <b>zipcode_98092</b> | -1.948e+04 | 7070.251 | -2.755 | 0.006 | -3.33e+04 | -5617.113 |
| <b>zipcode_98102</b> | 4.111e+05  | 1.05e+04 | 39.314 | 0.000 | 3.91e+05  | 4.32e+05  |
| <b>zipcode_98103</b> | 3.262e+05  | 6333.106 | 51.503 | 0.000 | 3.14e+05  | 3.39e+05  |
| <b>zipcode_98105</b> | 3.761e+05  | 8236.550 | 45.658 | 0.000 | 3.6e+05   | 3.92e+05  |
| <b>zipcode_98106</b> | 1.221e+05  | 6771.130 | 18.040 | 0.000 | 1.09e+05  | 1.35e+05  |
| <b>zipcode_98107</b> | 3.236e+05  | 7348.915 | 44.036 | 0.000 | 3.09e+05  | 3.38e+05  |
| <b>zipcode_98108</b> | 1.224e+05  | 7909.632 | 15.476 | 0.000 | 1.07e+05  | 1.38e+05  |
| <b>zipcode_98109</b> | 4.129e+05  | 1.05e+04 | 39.385 | 0.000 | 3.92e+05  | 4.33e+05  |
| <b>zipcode_98112</b> | 4.343e+05  | 8413.952 | 51.614 | 0.000 | 4.18e+05  | 4.51e+05  |
| <b>zipcode_98115</b> | 3.22e+05   | 6227.641 | 51.704 | 0.000 | 3.1e+05   | 3.34e+05  |
| <b>zipcode_98116</b> | 2.972e+05  | 7015.178 | 42.366 | 0.000 | 2.83e+05  | 3.11e+05  |
| <b>zipcode_98117</b> | 3.129e+05  | 6336.595 | 49.383 | 0.000 | 3e+05     | 3.25e+05  |
| <b>zipcode_98118</b> | 1.682e+05  | 6318.881 | 26.620 | 0.000 | 1.56e+05  | 1.81e+05  |
| <b>zipcode_98119</b> | 4.059e+05  | 8612.640 | 47.123 | 0.000 | 3.89e+05  | 4.23e+05  |
| <b>zipcode_98122</b> | 3.098e+05  | 7390.263 | 41.914 | 0.000 | 2.95e+05  | 3.24e+05  |
| <b>zipcode_98125</b> | 1.95e+05   | 6465.387 | 30.163 | 0.000 | 1.82e+05  | 2.08e+05  |
| <b>zipcode_98126</b> | 1.92e+05   | 6757.688 | 28.408 | 0.000 | 1.79e+05  | 2.05e+05  |
| <b>zipcode_98133</b> | 1.544e+05  | 6158.422 | 25.075 | 0.000 | 1.42e+05  | 1.66e+05  |
| <b>zipcode_98136</b> | 2.559e+05  | 7318.534 | 34.961 | 0.000 | 2.42e+05  | 2.7e+05   |

|                                       |            |          |        |       |           |           |
|---------------------------------------|------------|----------|--------|-------|-----------|-----------|
| <b>zipcode_98144</b>                  | 2.445e+05  | 7016.514 | 34.846 | 0.000 | 2.31e+05  | 2.58e+05  |
| <b>zipcode_98146</b>                  | 1.062e+05  | 7089.385 | 14.980 | 0.000 | 9.23e+04  | 1.2e+05   |
| <b>zipcode_98148</b>                  | 4.658e+04  | 1.22e+04 | 3.826  | 0.000 | 2.27e+04  | 7.04e+04  |
| <b>zipcode_98155</b>                  | 1.409e+05  | 6318.872 | 22.302 | 0.000 | 1.29e+05  | 1.53e+05  |
| <b>zipcode_98166</b>                  | 8.672e+04  | 7622.316 | 11.378 | 0.000 | 7.18e+04  | 1.02e+05  |
| <b>zipcode_98168</b>                  | 5.135e+04  | 7361.707 | 6.975  | 0.000 | 3.69e+04  | 6.58e+04  |
| <b>zipcode_98177</b>                  | 2.053e+05  | 7624.774 | 26.920 | 0.000 | 1.9e+05   | 2.2e+05   |
| <b>zipcode_98178</b>                  | 5.47e+04   | 7156.946 | 7.643  | 0.000 | 4.07e+04  | 6.87e+04  |
| <b>zipcode_98188</b>                  | 3.776e+04  | 8882.117 | 4.251  | 0.000 | 2.03e+04  | 5.52e+04  |
| <b>zipcode_98198</b>                  | 2.079e+04  | 7116.150 | 2.922  | 0.003 | 6844.626  | 3.47e+04  |
| <b>zipcode_98199</b>                  | 3.638e+05  | 7274.106 | 50.009 | 0.000 | 3.5e+05   | 3.78e+05  |
| <b>has_larger_sqft_than_neighbors</b> | -1.728e+04 | 1801.049 | -9.595 | 0.000 | -2.08e+04 | -1.38e+04 |

**Omnibus:** 1494.930    **Durbin-Watson:** 1.993

**Prob(Omnibus):** 0.000    **Jarque-Bera (JB):** 5136.551

**Skew:** 0.438    **Prob(JB):** 0.00

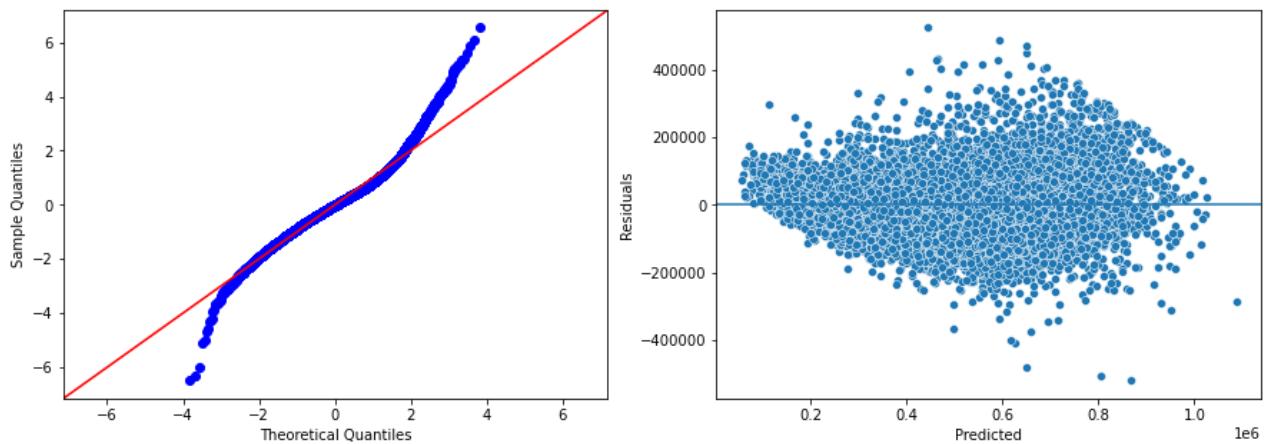
**Kurtosis:** 5.585    **Cond. No.** 9.83e+05

#### Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

[2] The condition number is large, 9.83e+05. This might indicate that there are

strong multicollinearity or other numerical problems.



After removing the lat and long columns due to high p-values, our model's R-squared values haven't changed but once again we are seeing that the residuals are better fitting the normality and homoscedasticity assumptions thanks to the IQR outlier removal process.

In [61]: df\_ohe

Out[61]: price bedrooms bathrooms sqft\_lot floors waterfront view condition grade sqft\_abov

|       | price    | bedrooms | bathrooms | sqft_lot | floors | waterfront | view | condition | grade | sqft_abov |
|-------|----------|----------|-----------|----------|--------|------------|------|-----------|-------|-----------|
| 0     | 221900.0 | 3        | 1.00      | 5650     | 1.0    | 0.0        | 0.0  | 3         | 7     | 118       |
| 1     | 538000.0 | 3        | 2.25      | 7242     | 2.0    | 0.0        | 0.0  | 3         | 7     | 217       |
| 2     | 180000.0 | 2        | 1.00      | 10000    | 1.0    | 0.0        | 0.0  | 3         | 6     | 77        |
| 3     | 604000.0 | 4        | 3.00      | 5000     | 1.0    | 0.0        | 0.0  | 5         | 7     | 105       |
| 4     | 510000.0 | 3        | 2.00      | 8080     | 1.0    | 0.0        | 0.0  | 3         | 8     | 168       |
| ...   | ...      | ...      | ...       | ...      | ...    | ...        | ...  | ...       | ...   | ...       |
| 21592 | 360000.0 | 3        | 2.50      | 1131     | 3.0    | 0.0        | 0.0  | 3         | 8     | 153       |
| 21593 | 400000.0 | 4        | 2.50      | 5813     | 2.0    | 0.0        | 0.0  | 3         | 8     | 231       |
| 21594 | 402101.0 | 2        | 0.75      | 1350     | 2.0    | 0.0        | 0.0  | 3         | 7     | 102       |
| 21595 | 400000.0 | 3        | 2.50      | 2388     | 2.0    | 0.0        | 0.0  | 3         | 8     | 160       |
| 21596 | 325000.0 | 2        | 0.75      | 1076     | 2.0    | 0.0        | 0.0  | 3         | 7     | 102       |

16556 rows × 83 columns

Since we removed outliers from all numeric columns including our target 'price', we are left with 16556 rows compared to the 21597 we started off with. It makes more sense to remove outliers based on the prices of the homes rather than removing outliers in every single column. This has a potential upside of allowing for there to be more data points and therefore a more accurate model overall.

## Removing Outliers Based on Price Only

```
In [62]: df_IQR_price = df_IQR_price[find_outliers_IQR(df_IQR_price['price'])==False]
df_IQR_price
```

|       | price    | bedrooms | bathrooms | sqft_lot | floors | waterfront | view | condition | grade | sqft_abov |
|-------|----------|----------|-----------|----------|--------|------------|------|-----------|-------|-----------|
| 0     | 221900.0 | 3        | 1.00      | 5650     | 1.0    | 0.0        | 0.0  | 3         | 7     | 118       |
| 1     | 538000.0 | 3        | 2.25      | 7242     | 2.0    | 0.0        | 0.0  | 3         | 7     | 217       |
| 2     | 180000.0 | 2        | 1.00      | 10000    | 1.0    | 0.0        | 0.0  | 3         | 6     | 77        |
| 3     | 604000.0 | 4        | 3.00      | 5000     | 1.0    | 0.0        | 0.0  | 5         | 7     | 105       |
| 4     | 510000.0 | 3        | 2.00      | 8080     | 1.0    | 0.0        | 0.0  | 3         | 8     | 168       |
| ...   | ...      | ...      | ...       | ...      | ...    | ...        | ...  | ...       | ...   | ...       |
| 21592 | 360000.0 | 3        | 2.50      | 1131     | 3.0    | 0.0        | 0.0  | 3         | 8     | 153       |
| 21593 | 400000.0 | 4        | 2.50      | 5813     | 2.0    | 0.0        | 0.0  | 3         | 8     | 231       |
| 21594 | 402101.0 | 2        | 0.75      | 1350     | 2.0    | 0.0        | 0.0  | 3         | 7     | 102       |
| 21595 | 400000.0 | 3        | 2.50      | 2388     | 2.0    | 0.0        | 0.0  | 3         | 8     | 160       |
| 21596 | 325000.0 | 2        | 0.75      | 1076     | 2.0    | 0.0        | 0.0  | 3         | 7     | 102       |

20439 rows × 85 columns

As seen above, we are left with approximately 4,000 more data points when we only remove outliers based on price. We still should take a look at the model and whether the residuals have adjusted similarly to the prior model.

## Model 2.2 - Addressing Outliers

In [63]: `model = model_lin_reg(df=df_IQR_price)`

| OLS Regression Results |                  |                     |             |       |           |           |        |  |  |  |  |
|------------------------|------------------|---------------------|-------------|-------|-----------|-----------|--------|--|--|--|--|
| Dep. Variable:         | price            | R-squared:          | 0.826       |       |           |           |        |  |  |  |  |
| Model:                 | OLS              | Adj. R-squared:     | 0.826       |       |           |           |        |  |  |  |  |
| Method:                | Least Squares    | F-statistic:        | 1152.       |       |           |           |        |  |  |  |  |
| Date:                  | Sat, 01 May 2021 | Prob (F-statistic): | 0.00        |       |           |           |        |  |  |  |  |
| Time:                  | 16:50:38         | Log-Likelihood:     | -2.6138e+05 |       |           |           |        |  |  |  |  |
| No. Observations:      | 20439            | AIC:                | 5.229e+05   |       |           |           |        |  |  |  |  |
| Df Residuals:          | 20354            | BIC:                | 5.236e+05   |       |           |           |        |  |  |  |  |
| Df Model:              | 84               |                     |             |       |           |           |        |  |  |  |  |
| Covariance Type:       | nonrobust        |                     |             |       |           |           |        |  |  |  |  |
|                        |                  | coef                | std err     | t     | P> t      | [0.025    | 0.975] |  |  |  |  |
| <b>Intercept</b>       | -1.461e+07       | 3.39e+06            | -4.308      | 0.000 | -2.13e+07 | -7.96e+06 |        |  |  |  |  |
| <b>bedrooms</b>        | -244.3677        | 850.207             | -0.287      | 0.774 | -1910.842 | 1422.106  |        |  |  |  |  |
| <b>bathrooms</b>       | 1.95e+04         | 1472.657            | 13.242      | 0.000 | 1.66e+04  | 2.24e+04  |        |  |  |  |  |
| <b>sqft_lot</b>        | 0.3002           | 0.017               | 17.975      | 0.000 | 0.268     | 0.333     |        |  |  |  |  |
| <b>floors</b>          | -2.479e+04       | 1761.730            | -14.073     | 0.000 | -2.82e+04 | -2.13e+04 |        |  |  |  |  |
| <b>waterfront</b>      | 1.501e+05        | 1.31e+04            | 11.488      | 0.000 | 1.24e+05  | 1.76e+05  |        |  |  |  |  |
| <b>view</b>            | 3.731e+04        | 1052.530            | 35.448      | 0.000 | 3.52e+04  | 3.94e+04  |        |  |  |  |  |
| <b>condition</b>       | 2.5e+04          | 1068.482            | 23.400      | 0.000 | 2.29e+04  | 2.71e+04  |        |  |  |  |  |
| <b>grade</b>           | 4.831e+04        | 1030.415            | 46.879      | 0.000 | 4.63e+04  | 5.03e+04  |        |  |  |  |  |
| <b>sqft_above</b>      | 131.7060         | 1.883               | 69.954      | 0.000 | 128.016   | 135.396   |        |  |  |  |  |
| <b>yr_built</b>        | -572.7285        | 36.232              | -15.807     | 0.000 | -643.747  | -501.710  |        |  |  |  |  |
| <b>lat</b>             | 1.545e+05        | 3.51e+04            | 4.402       | 0.000 | 8.57e+04  | 2.23e+05  |        |  |  |  |  |
| <b>long</b>            | -6.569e+04       | 2.51e+04            | -2.621      | 0.009 | -1.15e+05 | -1.66e+04 |        |  |  |  |  |
| <b>renovated</b>       | 3.36e+04         | 3726.684            | 9.017       | 0.000 | 2.63e+04  | 4.09e+04  |        |  |  |  |  |
| <b>has_basement</b>    | 4.727e+04        | 1670.349            | 28.300      | 0.000 | 4.4e+04   | 5.05e+04  |        |  |  |  |  |
| <b>zipcode_98002</b>   | 6083.3919        | 7819.865            | 0.778       | 0.437 | -9244.173 | 2.14e+04  |        |  |  |  |  |

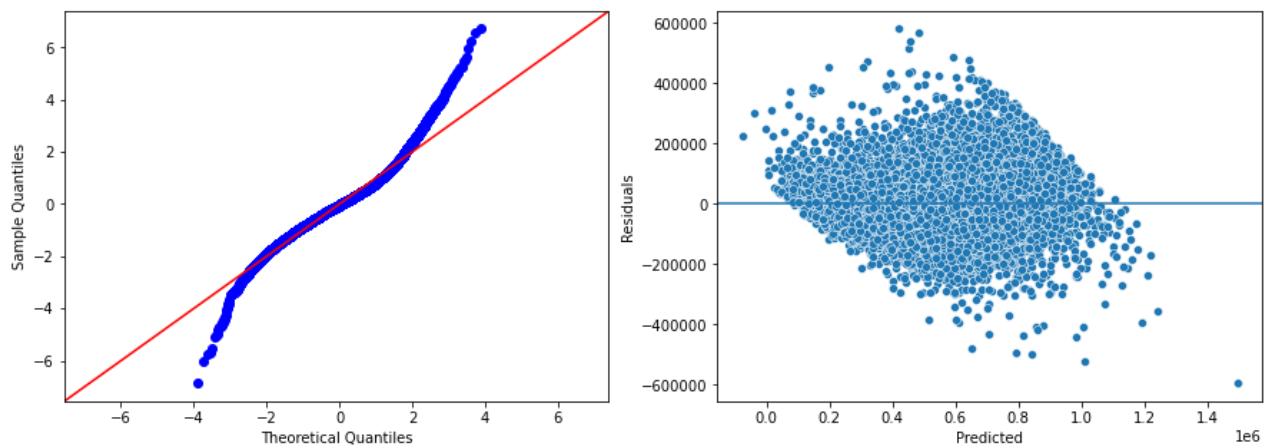
|                      | final_notebook-Copy1 |          |        |       |           |           |  |
|----------------------|----------------------|----------|--------|-------|-----------|-----------|--|
| <b>zipcode_98003</b> | -1.129e+04           | 6991.183 | -1.615 | 0.106 | -2.5e+04  | 2411.862  |  |
| <b>zipcode_98004</b> | 4.749e+05            | 1.37e+04 | 34.573 | 0.000 | 4.48e+05  | 5.02e+05  |  |
| <b>zipcode_98005</b> | 2.94e+05             | 1.39e+04 | 21.181 | 0.000 | 2.67e+05  | 3.21e+05  |  |
| <b>zipcode_98006</b> | 2.484e+05            | 1.14e+04 | 21.769 | 0.000 | 2.26e+05  | 2.71e+05  |  |
| <b>zipcode_98007</b> | 2.177e+05            | 1.42e+04 | 15.285 | 0.000 | 1.9e+05   | 2.46e+05  |  |
| <b>zipcode_98008</b> | 2.006e+05            | 1.36e+04 | 14.753 | 0.000 | 1.74e+05  | 2.27e+05  |  |
| <b>zipcode_98010</b> | 9.965e+04            | 1.2e+04  | 8.310  | 0.000 | 7.61e+04  | 1.23e+05  |  |
| <b>zipcode_98011</b> | 8.511e+04            | 1.76e+04 | 4.825  | 0.000 | 5.05e+04  | 1.2e+05   |  |
| <b>zipcode_98014</b> | 8.339e+04            | 1.94e+04 | 4.297  | 0.000 | 4.54e+04  | 1.21e+05  |  |
| <b>zipcode_98019</b> | 5.97e+04             | 1.91e+04 | 3.127  | 0.002 | 2.23e+04  | 9.71e+04  |  |
| <b>zipcode_98022</b> | 2.401e+04            | 1.05e+04 | 2.291  | 0.022 | 3470.958  | 4.45e+04  |  |
| <b>zipcode_98023</b> | -2.911e+04           | 6450.507 | -4.513 | 0.000 | -4.18e+04 | -1.65e+04 |  |
| <b>zipcode_98024</b> | 1.335e+05            | 1.72e+04 | 7.770  | 0.000 | 9.98e+04  | 1.67e+05  |  |
| <b>zipcode_98027</b> | 1.777e+05            | 1.16e+04 | 15.365 | 0.000 | 1.55e+05  | 2e+05     |  |
| <b>zipcode_98028</b> | 7.243e+04            | 1.72e+04 | 4.221  | 0.000 | 3.88e+04  | 1.06e+05  |  |
| <b>zipcode_98029</b> | 2.089e+05            | 1.32e+04 | 15.786 | 0.000 | 1.83e+05  | 2.35e+05  |  |
| <b>zipcode_98030</b> | 1427.2186            | 7704.548 | 0.185  | 0.853 | -1.37e+04 | 1.65e+04  |  |
| <b>zipcode_98031</b> | 2295.0973            | 8044.248 | 0.285  | 0.775 | -1.35e+04 | 1.81e+04  |  |
| <b>zipcode_98032</b> | -1.8e+04             | 9302.376 | -1.935 | 0.053 | -3.62e+04 | 235.348   |  |
| <b>zipcode_98033</b> | 2.702e+05            | 1.48e+04 | 18.204 | 0.000 | 2.41e+05  | 2.99e+05  |  |
| <b>zipcode_98034</b> | 1.25e+05             | 1.58e+04 | 7.906  | 0.000 | 9.4e+04   | 1.56e+05  |  |
| <b>zipcode_98038</b> | 5e+04                | 8718.347 | 5.735  | 0.000 | 3.29e+04  | 6.71e+04  |  |
| <b>zipcode_98039</b> | 6.029e+05            | 3.75e+04 | 16.072 | 0.000 | 5.29e+05  | 6.76e+05  |  |
| <b>zipcode_98040</b> | 3.904e+05            | 1.22e+04 | 32.021 | 0.000 | 3.67e+05  | 4.14e+05  |  |
| <b>zipcode_98042</b> | 1.043e+04            | 7403.975 | 1.409  | 0.159 | -4081.403 | 2.49e+04  |  |
| <b>zipcode_98045</b> | 1.112e+05            | 1.62e+04 | 6.866  | 0.000 | 7.94e+04  | 1.43e+05  |  |
| <b>zipcode_98052</b> | 2.092e+05            | 1.5e+04  | 13.929 | 0.000 | 1.8e+05   | 2.39e+05  |  |
| <b>zipcode_98053</b> | 2.011e+05            | 1.61e+04 | 12.466 | 0.000 | 1.7e+05   | 2.33e+05  |  |
| <b>zipcode_98055</b> | 2.19e+04             | 8989.534 | 2.436  | 0.015 | 4281.690  | 3.95e+04  |  |
| <b>zipcode_98056</b> | 7.814e+04            | 9814.522 | 7.961  | 0.000 | 5.89e+04  | 9.74e+04  |  |
| <b>zipcode_98058</b> | 2.414e+04            | 8506.076 | 2.838  | 0.005 | 7466.703  | 4.08e+04  |  |
| <b>zipcode_98059</b> | 8.623e+04            | 9615.356 | 8.968  | 0.000 | 6.74e+04  | 1.05e+05  |  |
| <b>zipcode_98065</b> | 1.264e+05            | 1.49e+04 | 8.473  | 0.000 | 9.72e+04  | 1.56e+05  |  |
| <b>zipcode_98070</b> | 7.2e+04              | 1.14e+04 | 6.334  | 0.000 | 4.97e+04  | 9.43e+04  |  |
| <b>zipcode_98072</b> | 1.227e+05            | 1.76e+04 | 6.978  | 0.000 | 8.82e+04  | 1.57e+05  |  |

|                                       |            |          |        |       |           |           |
|---------------------------------------|------------|----------|--------|-------|-----------|-----------|
| <b>zipcode_98074</b>                  | 1.88e+05   | 1.42e+04 | 13.231 | 0.000 | 1.6e+05   | 2.16e+05  |
| <b>zipcode_98075</b>                  | 2.145e+05  | 1.37e+04 | 15.671 | 0.000 | 1.88e+05  | 2.41e+05  |
| <b>zipcode_98077</b>                  | 1.269e+05  | 1.83e+04 | 6.929  | 0.000 | 9.1e+04   | 1.63e+05  |
| <b>zipcode_98092</b>                  | -7663.2632 | 7003.594 | -1.094 | 0.274 | -2.14e+04 | 6064.345  |
| <b>zipcode_98102</b>                  | 3.42e+05   | 1.55e+04 | 22.013 | 0.000 | 3.12e+05  | 3.72e+05  |
| <b>zipcode_98103</b>                  | 2.507e+05  | 1.43e+04 | 17.591 | 0.000 | 2.23e+05  | 2.79e+05  |
| <b>zipcode_98105</b>                  | 3.089e+05  | 1.49e+04 | 20.695 | 0.000 | 2.8e+05   | 3.38e+05  |
| <b>zipcode_98106</b>                  | 7.308e+04  | 1.05e+04 | 6.969  | 0.000 | 5.25e+04  | 9.36e+04  |
| <b>zipcode_98107</b>                  | 2.481e+05  | 1.47e+04 | 16.910 | 0.000 | 2.19e+05  | 2.77e+05  |
| <b>zipcode_98108</b>                  | 7.486e+04  | 1.16e+04 | 6.476  | 0.000 | 5.22e+04  | 9.75e+04  |
| <b>zipcode_98109</b>                  | 3.432e+05  | 1.56e+04 | 21.980 | 0.000 | 3.13e+05  | 3.74e+05  |
| <b>zipcode_98112</b>                  | 3.687e+05  | 1.4e+04  | 26.416 | 0.000 | 3.41e+05  | 3.96e+05  |
| <b>zipcode_98115</b>                  | 2.535e+05  | 1.45e+04 | 17.488 | 0.000 | 2.25e+05  | 2.82e+05  |
| <b>zipcode_98116</b>                  | 2.366e+05  | 1.18e+04 | 20.095 | 0.000 | 2.13e+05  | 2.6e+05   |
| <b>zipcode_98117</b>                  | 2.387e+05  | 1.47e+04 | 16.278 | 0.000 | 2.1e+05   | 2.67e+05  |
| <b>zipcode_98118</b>                  | 1.236e+05  | 1.02e+04 | 12.064 | 0.000 | 1.04e+05  | 1.44e+05  |
| <b>zipcode_98119</b>                  | 3.321e+05  | 1.45e+04 | 22.835 | 0.000 | 3.04e+05  | 3.61e+05  |
| <b>zipcode_98122</b>                  | 2.458e+05  | 1.28e+04 | 19.265 | 0.000 | 2.21e+05  | 2.71e+05  |
| <b>zipcode_98125</b>                  | 1.291e+05  | 1.56e+04 | 8.261  | 0.000 | 9.85e+04  | 1.6e+05   |
| <b>zipcode_98126</b>                  | 1.385e+05  | 1.08e+04 | 12.873 | 0.000 | 1.17e+05  | 1.6e+05   |
| <b>zipcode_98133</b>                  | 7.66e+04   | 1.61e+04 | 4.744  | 0.000 | 4.5e+04   | 1.08e+05  |
| <b>zipcode_98136</b>                  | 2.047e+05  | 1.11e+04 | 18.520 | 0.000 | 1.83e+05  | 2.26e+05  |
| <b>zipcode_98144</b>                  | 1.918e+05  | 1.19e+04 | 16.157 | 0.000 | 1.68e+05  | 2.15e+05  |
| <b>zipcode_98146</b>                  | 7.139e+04  | 9838.756 | 7.256  | 0.000 | 5.21e+04  | 9.07e+04  |
| <b>zipcode_98148</b>                  | 3.006e+04  | 1.32e+04 | 2.273  | 0.023 | 4139.841  | 5.6e+04   |
| <b>zipcode_98155</b>                  | 6.584e+04  | 1.68e+04 | 3.919  | 0.000 | 3.29e+04  | 9.88e+04  |
| <b>zipcode_98166</b>                  | 7.007e+04  | 9010.639 | 7.777  | 0.000 | 5.24e+04  | 8.77e+04  |
| <b>zipcode_98168</b>                  | 1.879e+04  | 9482.888 | 1.981  | 0.048 | 201.759   | 3.74e+04  |
| <b>zipcode_98177</b>                  | 1.401e+05  | 1.69e+04 | 8.267  | 0.000 | 1.07e+05  | 1.73e+05  |
| <b>zipcode_98178</b>                  | 2.552e+04  | 9804.820 | 2.602  | 0.009 | 6298.076  | 4.47e+04  |
| <b>zipcode_98188</b>                  | 1.038e+04  | 9994.135 | 1.039  | 0.299 | -9209.027 | 3e+04     |
| <b>zipcode_98198</b>                  | 5223.4014  | 7582.058 | 0.689  | 0.491 | -9638.043 | 2.01e+04  |
| <b>zipcode_98199</b>                  | 2.926e+05  | 1.4e+04  | 20.843 | 0.000 | 2.65e+05  | 3.2e+05   |
| <b>has_larger_sqft_than_neighbors</b> | -1.718e+04 | 1719.318 | -9.993 | 0.000 | -2.06e+04 | -1.38e+04 |

|                       |          |                          |          |
|-----------------------|----------|--------------------------|----------|
| <b>Omnibus:</b>       | 1997.283 | <b>Durbin-Watson:</b>    | 1.988    |
| <b>Prob(Omnibus):</b> | 0.000    | <b>Jarque-Bera (JB):</b> | 7373.063 |
| <b>Skew:</b>          | 0.457    | <b>Prob(JB):</b>         | 0.00     |
| <b>Kurtosis:</b>      | 5.797    | <b>Cond. No.</b>         | 2.38e+08 |

Notes:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [2] The condition number is large, 2.38e+08. This might indicate that there are strong multicollinearity or other numerical problems.



```
In [64]: model.pvalues[model.pvalues>0.05]
```

```
Out[64]: bedrooms      0.773793
zipcode_98002    0.436613
zipcode_98003    0.106306
zipcode_98030    0.853040
zipcode_98031    0.775410
zipcode_98032    0.053031
zipcode_98042    0.158899
zipcode_98092    0.273885
zipcode_98188    0.298986
zipcode_98198    0.490884
dtype: float64
```

Our model has higher R-squared and adjusted R-squared values at 0.826 compared to the prior model's 0.820 and 0.819. It is important to note that we achieved a better fit for our model based on these numbers and kept more data points compared to model 2.1. We will move forward with this model instead of using the prior one due to this reason. Additionally, compared to model 2.1, the 'lat' and 'long' columns are significant in this model so we are keeping them in.

Interestingly, we are seeing that bedrooms are no longer a significant coefficient, and therefore not a significant parameter to define the sales price of a home based on our model. This could be due to the overall above ground square footage being a better predictor of a home's value rather than bedroom counts.

```
In [65]: df_IQR_price.drop('bedrooms', axis=1, inplace=True)
model_non_scaled = model_lin_reg(df=df_IQR_price)
```

C:\Users\berke\anaconda3\envs\learn-env\lib\site-packages\pandas\core\frame.py:4163: SettingWithCopyWarning:

A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

```
    return super().drop(
```

OLS Regression Results

|                          |                  |                            |             |
|--------------------------|------------------|----------------------------|-------------|
| <b>Dep. Variable:</b>    | price            | <b>R-squared:</b>          | 0.826       |
| <b>Model:</b>            | OLS              | <b>Adj. R-squared:</b>     | 0.826       |
| <b>Method:</b>           | Least Squares    | <b>F-statistic:</b>        | 1166.       |
| <b>Date:</b>             | Sat, 01 May 2021 | <b>Prob (F-statistic):</b> | 0.00        |
| <b>Time:</b>             | 16:50:39         | <b>Log-Likelihood:</b>     | -2.6138e+05 |
| <b>No. Observations:</b> | 20439            | <b>AIC:</b>                | 5.229e+05   |
| <b>Df Residuals:</b>     | 20355            | <b>BIC:</b>                | 5.236e+05   |
| <b>Df Model:</b>         | 83               |                            |             |
| <b>Covariance Type:</b>  | nonrobust        |                            |             |

|                      |            | coef     | std err | t     | P> t      | [0.025    | 0.975] |
|----------------------|------------|----------|---------|-------|-----------|-----------|--------|
| <b>Intercept</b>     | -1.461e+07 | 3.39e+06 | -4.308  | 0.000 | -2.13e+07 | -7.96e+06 |        |
| <b>bathrooms</b>     | 1.94e+04   | 1427.215 | 13.590  | 0.000 | 1.66e+04  | 2.22e+04  |        |
| <b>sqft_lot</b>      | 0.3005     | 0.017    | 18.011  | 0.000 | 0.268     | 0.333     |        |
| <b>floors</b>        | -2.477e+04 | 1759.383 | -14.078 | 0.000 | -2.82e+04 | -2.13e+04 |        |
| <b>waterfront</b>    | 1.501e+05  | 1.31e+04 | 11.495  | 0.000 | 1.25e+05  | 1.76e+05  |        |
| <b>view</b>          | 3.732e+04  | 1051.829 | 35.482  | 0.000 | 3.53e+04  | 3.94e+04  |        |
| <b>condition</b>     | 2.499e+04  | 1067.722 | 23.406  | 0.000 | 2.29e+04  | 2.71e+04  |        |
| <b>grade</b>         | 4.833e+04  | 1025.815 | 47.117  | 0.000 | 4.63e+04  | 5.03e+04  |        |
| <b>sqft_above</b>    | 131.5554   | 1.808    | 72.748  | 0.000 | 128.011   | 135.100   |        |
| <b>yr_built</b>      | -571.7167  | 36.060   | -15.855 | 0.000 | -642.397  | -501.036  |        |
| <b>lat</b>           | 1.546e+05  | 3.51e+04 | 4.406   | 0.000 | 8.58e+04  | 2.23e+05  |        |
| <b>long</b>          | -6.562e+04 | 2.51e+04 | -2.618  | 0.009 | -1.15e+05 | -1.65e+04 |        |
| <b>renovated</b>     | 3.364e+04  | 3724.476 | 9.032   | 0.000 | 2.63e+04  | 4.09e+04  |        |
| <b>has_basement</b>  | 4.717e+04  | 1636.087 | 28.833  | 0.000 | 4.4e+04   | 5.04e+04  |        |
| <b>zipcode_98002</b> | 6079.6824  | 7819.678 | 0.777   | 0.437 | -9247.516 | 2.14e+04  |        |
| <b>zipcode_98003</b> | -1.128e+04 | 6990.820 | -1.613  | 0.107 | -2.5e+04  | 2426.563  |        |
| <b>zipcode_98004</b> | 4.749e+05  | 1.37e+04 | 34.573  | 0.000 | 4.48e+05  | 5.02e+05  |        |
| <b>zipcode_98005</b> | 2.939e+05  | 1.39e+04 | 21.180  | 0.000 | 2.67e+05  | 3.21e+05  |        |
| <b>zipcode_98006</b> | 2.484e+05  | 1.14e+04 | 21.768  | 0.000 | 2.26e+05  | 2.71e+05  |        |
| <b>zipcode_98007</b> | 2.176e+05  | 1.42e+04 | 15.283  | 0.000 | 1.9e+05   | 2.46e+05  |        |

|                      |            |          |        |       |           |           |
|----------------------|------------|----------|--------|-------|-----------|-----------|
| <b>zipcode_98008</b> | 2.005e+05  | 1.36e+04 | 14.751 | 0.000 | 1.74e+05  | 2.27e+05  |
| <b>zipcode_98010</b> | 9.968e+04  | 1.2e+04  | 8.313  | 0.000 | 7.62e+04  | 1.23e+05  |
| <b>zipcode_98011</b> | 8.508e+04  | 1.76e+04 | 4.824  | 0.000 | 5.05e+04  | 1.2e+05   |
| <b>zipcode_98014</b> | 8.342e+04  | 1.94e+04 | 4.299  | 0.000 | 4.54e+04  | 1.21e+05  |
| <b>zipcode_98019</b> | 5.968e+04  | 1.91e+04 | 3.126  | 0.002 | 2.23e+04  | 9.71e+04  |
| <b>zipcode_98022</b> | 2.404e+04  | 1.05e+04 | 2.294  | 0.022 | 3502.443  | 4.46e+04  |
| <b>zipcode_98023</b> | -2.91e+04  | 6450.225 | -4.512 | 0.000 | -4.17e+04 | -1.65e+04 |
| <b>zipcode_98024</b> | 1.335e+05  | 1.72e+04 | 7.771  | 0.000 | 9.99e+04  | 1.67e+05  |
| <b>zipcode_98027</b> | 1.777e+05  | 1.16e+04 | 15.366 | 0.000 | 1.55e+05  | 2e+05     |
| <b>zipcode_98028</b> | 7.239e+04  | 1.72e+04 | 4.219  | 0.000 | 3.88e+04  | 1.06e+05  |
| <b>zipcode_98029</b> | 2.089e+05  | 1.32e+04 | 15.786 | 0.000 | 1.83e+05  | 2.35e+05  |
| <b>zipcode_98030</b> | 1409.0238  | 7704.114 | 0.183  | 0.855 | -1.37e+04 | 1.65e+04  |
| <b>zipcode_98031</b> | 2269.2724  | 8043.564 | 0.282  | 0.778 | -1.35e+04 | 1.8e+04   |
| <b>zipcode_98032</b> | -1.804e+04 | 9301.251 | -1.939 | 0.053 | -3.63e+04 | 195.624   |
| <b>zipcode_98033</b> | 2.701e+05  | 1.48e+04 | 18.203 | 0.000 | 2.41e+05  | 2.99e+05  |
| <b>zipcode_98034</b> | 1.25e+05   | 1.58e+04 | 7.904  | 0.000 | 9.4e+04   | 1.56e+05  |
| <b>zipcode_98038</b> | 5.001e+04  | 8718.135 | 5.736  | 0.000 | 3.29e+04  | 6.71e+04  |
| <b>zipcode_98039</b> | 6.029e+05  | 3.75e+04 | 16.071 | 0.000 | 5.29e+05  | 6.76e+05  |
| <b>zipcode_98040</b> | 3.904e+05  | 1.22e+04 | 32.021 | 0.000 | 3.66e+05  | 4.14e+05  |
| <b>zipcode_98042</b> | 1.043e+04  | 7403.795 | 1.408  | 0.159 | -4085.059 | 2.49e+04  |
| <b>zipcode_98045</b> | 1.111e+05  | 1.62e+04 | 6.865  | 0.000 | 7.94e+04  | 1.43e+05  |
| <b>zipcode_98052</b> | 2.092e+05  | 1.5e+04  | 13.927 | 0.000 | 1.8e+05   | 2.39e+05  |
| <b>zipcode_98053</b> | 2.012e+05  | 1.61e+04 | 12.472 | 0.000 | 1.7e+05   | 2.33e+05  |
| <b>zipcode_98055</b> | 2.192e+04  | 8989.182 | 2.438  | 0.015 | 4297.280  | 3.95e+04  |
| <b>zipcode_98056</b> | 7.812e+04  | 9814.117 | 7.960  | 0.000 | 5.89e+04  | 9.74e+04  |
| <b>zipcode_98058</b> | 2.411e+04  | 8505.148 | 2.834  | 0.005 | 7436.354  | 4.08e+04  |
| <b>zipcode_98059</b> | 8.62e+04   | 9614.416 | 8.965  | 0.000 | 6.74e+04  | 1.05e+05  |
| <b>zipcode_98065</b> | 1.264e+05  | 1.49e+04 | 8.474  | 0.000 | 9.72e+04  | 1.56e+05  |
| <b>zipcode_98070</b> | 7.211e+04  | 1.14e+04 | 6.348  | 0.000 | 4.98e+04  | 9.44e+04  |
| <b>zipcode_98072</b> | 1.227e+05  | 1.76e+04 | 6.978  | 0.000 | 8.82e+04  | 1.57e+05  |
| <b>zipcode_98074</b> | 1.88e+05   | 1.42e+04 | 13.230 | 0.000 | 1.6e+05   | 2.16e+05  |
| <b>zipcode_98075</b> | 2.145e+05  | 1.37e+04 | 15.671 | 0.000 | 1.88e+05  | 2.41e+05  |
| <b>zipcode_98077</b> | 1.269e+05  | 1.83e+04 | 6.929  | 0.000 | 9.1e+04   | 1.63e+05  |
| <b>zipcode_98092</b> | -7660.1869 | 7003.428 | -1.094 | 0.274 | -2.14e+04 | 6067.095  |
| <b>zipcode_98102</b> | 3.42e+05   | 1.55e+04 | 22.018 | 0.000 | 3.12e+05  | 3.72e+05  |

|                                       |            |          |         |       |           |           |
|---------------------------------------|------------|----------|---------|-------|-----------|-----------|
| <b>zipcode_98103</b>                  | 2.507e+05  | 1.43e+04 | 17.593  | 0.000 | 2.23e+05  | 2.79e+05  |
| <b>zipcode_98105</b>                  | 3.089e+05  | 1.49e+04 | 20.694  | 0.000 | 2.8e+05   | 3.38e+05  |
| <b>zipcode_98106</b>                  | 7.308e+04  | 1.05e+04 | 6.970   | 0.000 | 5.25e+04  | 9.36e+04  |
| <b>zipcode_98107</b>                  | 2.482e+05  | 1.47e+04 | 16.914  | 0.000 | 2.19e+05  | 2.77e+05  |
| <b>zipcode_98108</b>                  | 7.488e+04  | 1.16e+04 | 6.478   | 0.000 | 5.22e+04  | 9.75e+04  |
| <b>zipcode_98109</b>                  | 3.433e+05  | 1.56e+04 | 21.986  | 0.000 | 3.13e+05  | 3.74e+05  |
| <b>zipcode_98112</b>                  | 3.688e+05  | 1.4e+04  | 26.420  | 0.000 | 3.41e+05  | 3.96e+05  |
| <b>zipcode_98115</b>                  | 2.535e+05  | 1.45e+04 | 17.490  | 0.000 | 2.25e+05  | 2.82e+05  |
| <b>zipcode_98116</b>                  | 2.366e+05  | 1.18e+04 | 20.103  | 0.000 | 2.14e+05  | 2.6e+05   |
| <b>zipcode_98117</b>                  | 2.387e+05  | 1.47e+04 | 16.282  | 0.000 | 2.1e+05   | 2.67e+05  |
| <b>zipcode_98118</b>                  | 1.236e+05  | 1.02e+04 | 12.067  | 0.000 | 1.04e+05  | 1.44e+05  |
| <b>zipcode_98119</b>                  | 3.321e+05  | 1.45e+04 | 22.841  | 0.000 | 3.04e+05  | 3.61e+05  |
| <b>zipcode_98122</b>                  | 2.459e+05  | 1.28e+04 | 19.268  | 0.000 | 2.21e+05  | 2.71e+05  |
| <b>zipcode_98125</b>                  | 1.291e+05  | 1.56e+04 | 8.260   | 0.000 | 9.85e+04  | 1.6e+05   |
| <b>zipcode_98126</b>                  | 1.386e+05  | 1.08e+04 | 12.884  | 0.000 | 1.17e+05  | 1.6e+05   |
| <b>zipcode_98133</b>                  | 7.658e+04  | 1.61e+04 | 4.743   | 0.000 | 4.49e+04  | 1.08e+05  |
| <b>zipcode_98136</b>                  | 2.047e+05  | 1.1e+04  | 18.531  | 0.000 | 1.83e+05  | 2.26e+05  |
| <b>zipcode_98144</b>                  | 1.918e+05  | 1.19e+04 | 16.160  | 0.000 | 1.69e+05  | 2.15e+05  |
| <b>zipcode_98146</b>                  | 7.139e+04  | 9838.519 | 7.257   | 0.000 | 5.21e+04  | 9.07e+04  |
| <b>zipcode_98148</b>                  | 3.009e+04  | 1.32e+04 | 2.275   | 0.023 | 4163.796  | 5.6e+04   |
| <b>zipcode_98155</b>                  | 6.579e+04  | 1.68e+04 | 3.917   | 0.000 | 3.29e+04  | 9.87e+04  |
| <b>zipcode_98166</b>                  | 7.008e+04  | 9010.377 | 7.778   | 0.000 | 5.24e+04  | 8.77e+04  |
| <b>zipcode_98168</b>                  | 1.88e+04   | 9482.544 | 1.983   | 0.047 | 216.730   | 3.74e+04  |
| <b>zipcode_98177</b>                  | 1.401e+05  | 1.69e+04 | 8.267   | 0.000 | 1.07e+05  | 1.73e+05  |
| <b>zipcode_98178</b>                  | 2.549e+04  | 9804.095 | 2.600   | 0.009 | 6270.930  | 4.47e+04  |
| <b>zipcode_98188</b>                  | 1.035e+04  | 9993.181 | 1.035   | 0.301 | -9241.852 | 2.99e+04  |
| <b>zipcode_98198</b>                  | 5236.4108  | 7581.752 | 0.691   | 0.490 | -9624.433 | 2.01e+04  |
| <b>zipcode_98199</b>                  | 2.926e+05  | 1.4e+04  | 20.850  | 0.000 | 2.65e+05  | 3.2e+05   |
| <b>has_larger_sqft_than_neighbors</b> | -1.719e+04 | 1718.780 | -10.003 | 0.000 | -2.06e+04 | -1.38e+04 |

**Omnibus:** 1998.403    **Durbin-Watson:** 1.987

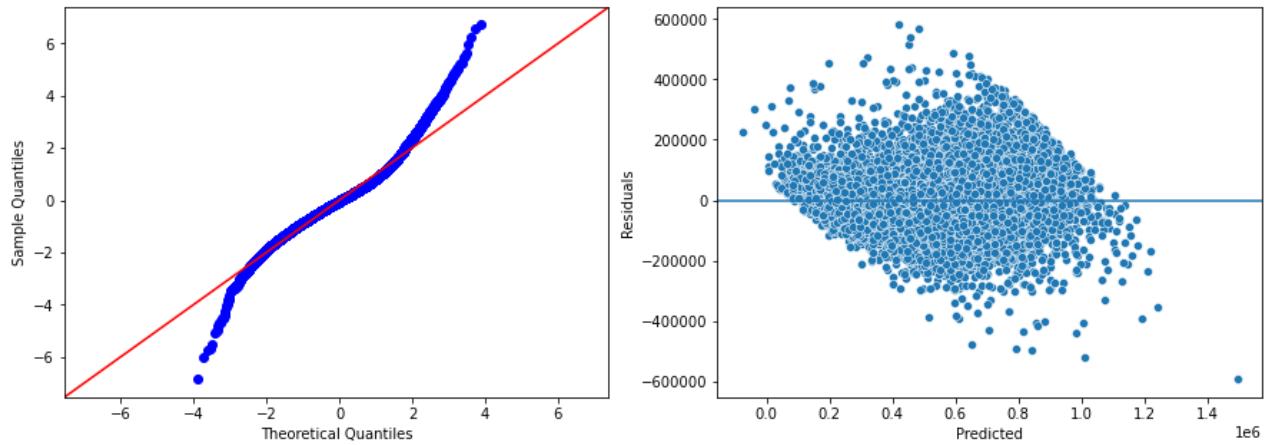
**Prob(Omnibus):** 0.000    **Jarque-Bera (JB):** 7377.212

**Skew:** 0.458    **Prob(JB):** 0.00

**Kurtosis:** 5.797    **Cond. No.** 2.38e+08

Notes:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [2] The condition number is large, 2.38e+08. This might indicate that there are strong multicollinearity or other numerical problems.



After removing the bedrooms column, we are left with no issues with any of the p-values for our parameters and as discussed above we still have R-squared and adjusted R-squared values of 0.826 with our residuals meeting the normality and homoscedasticity assumptions.

## Scaling

Now that we have an accurate model we can go ahead and scale the data to see which of the parameters affect the price the most relative to other parameters.

```
In [66]: from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
```

```
In [67]: df_IQR_price.columns
```

```
Out[67]: Index(['price', 'bathrooms', 'sqft_lot', 'floors', 'waterfront', 'view',
       'condition', 'grade', 'sqft_above', 'yr_built', 'lat', 'long',
       'renovated', 'has_basement', 'zipcode_98002', 'zipcode_98003',
       'zipcode_98004', 'zipcode_98005', 'zipcode_98006', 'zipcode_98007',
       'zipcode_98008', 'zipcode_98010', 'zipcode_98011', 'zipcode_98014',
       'zipcode_98019', 'zipcode_98022', 'zipcode_98023', 'zipcode_98024',
       'zipcode_98027', 'zipcode_98028', 'zipcode_98029', 'zipcode_98030',
       'zipcode_98031', 'zipcode_98032', 'zipcode_98033', 'zipcode_98034',
       'zipcode_98038', 'zipcode_98039', 'zipcode_98040', 'zipcode_98042',
       'zipcode_98045', 'zipcode_98052', 'zipcode_98053', 'zipcode_98055',
       'zipcode_98056', 'zipcode_98058', 'zipcode_98059', 'zipcode_98065',
       'zipcode_98070', 'zipcode_98072', 'zipcode_98074', 'zipcode_98075',
       'zipcode_98077', 'zipcode_98092', 'zipcode_98102', 'zipcode_98103',
       'zipcode_98105', 'zipcode_98106', 'zipcode_98107', 'zipcode_98108',
       'zipcode_98109', 'zipcode_98112', 'zipcode_98115', 'zipcode_98116',
       'zipcode_98117', 'zipcode_98118', 'zipcode_98119', 'zipcode_98122',
       'zipcode_98125', 'zipcode_98126', 'zipcode_98133', 'zipcode_98136',
       'zipcode_98144', 'zipcode_98146', 'zipcode_98148', 'zipcode_98155',
       'zipcode_98166', 'zipcode_98168', 'zipcode_98177', 'zipcode_98178',
       'zipcode_98188', 'zipcode_98198', 'zipcode_98199',
       'has_larger_sqft_than_neighbors'],
      dtype='object')
```

```
In [68]: numeric_cols = [col for col in df_IQR_price.columns if (col.startswith('zipcode'))==False]
```

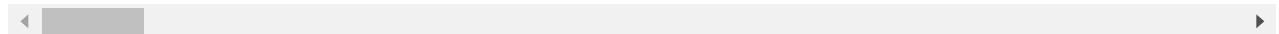
```
In [69]: numeric_cols
```

```
Out[69]: ['bathrooms',
 'sqft_lot',
 'floors',
 'waterfront',
 'view',
 'condition',
 'grade',
 'sqft_above',
 'yr_built',
 'lat',
 'long',
 'renovated']
```

```
In [70]: df_scaled = df_IQR_price.copy()
df_scaled[numeric_cols] = scaler.fit_transform(df_scaled[numeric_cols])
df_scaled.head()
```

|   | price    | bathrooms | sqft_lot  | floors    | waterfront | view      | condition | grade     | sqft_above |
|---|----------|-----------|-----------|-----------|------------|-----------|-----------|-----------|------------|
| 0 | 221900.0 | -1.481158 | -0.223340 | -0.886223 | -0.050015  | -0.268494 | -0.628492 | -0.512196 | -0.726541  |
| 1 | 538000.0 | 0.278964  | -0.183549 | 0.977307  | -0.050015  | -0.268494 | -0.628492 | -0.512196 | 0.635417   |
| 2 | 180000.0 | -1.481158 | -0.114614 | -0.886223 | -0.050015  | -0.268494 | -0.628492 | -1.477415 | -1.290585  |
| 3 | 604000.0 | 1.335036  | -0.239586 | -0.886223 | -0.050015  | -0.268494 | 2.462772  | -0.512196 | -0.905384  |
| 4 | 510000.0 | -0.073061 | -0.162603 | -0.886223 | -0.050015  | -0.268494 | -0.628492 | 0.453023  | -0.038683  |

5 rows × 84 columns



## Model 3 - Scaled Model

```
In [71]: model = model_lin_reg(df=df_scaled)
```

OLS Regression Results

|                          |                  |                            |             |
|--------------------------|------------------|----------------------------|-------------|
| <b>Dep. Variable:</b>    | price            | <b>R-squared:</b>          | 0.826       |
| <b>Model:</b>            | OLS              | <b>Adj. R-squared:</b>     | 0.826       |
| <b>Method:</b>           | Least Squares    | <b>F-statistic:</b>        | 1166.       |
| <b>Date:</b>             | Sat, 01 May 2021 | <b>Prob (F-statistic):</b> | 0.00        |
| <b>Time:</b>             | 16:50:41         | <b>Log-Likelihood:</b>     | -2.6138e+05 |
| <b>No. Observations:</b> | 20439            | <b>AIC:</b>                | 5.229e+05   |
| <b>Df Residuals:</b>     | 20355            | <b>BIC:</b>                | 5.236e+05   |
| <b>Df Model:</b>         | 83               |                            |             |
| <b>Covariance Type:</b>  | nonrobust        |                            |             |

|                  | coef      | std err | t      | P> t  | [0.025   | 0.975]  |
|------------------|-----------|---------|--------|-------|----------|---------|
| <b>Intercept</b> | 3.303e+05 | 1e+04   | 32.877 | 0.000 | 3.11e+05 | 3.5e+05 |

|                      |            |          |         |       |           |           |
|----------------------|------------|----------|---------|-------|-----------|-----------|
| <b>bathrooms</b>     | 1.377e+04  | 1013.577 | 13.590  | 0.000 | 1.18e+04  | 1.58e+04  |
| <b>sqft_lot</b>      | 1.202e+04  | 667.475  | 18.011  | 0.000 | 1.07e+04  | 1.33e+04  |
| <b>floors</b>        | -1.329e+04 | 944.113  | -14.078 | 0.000 | -1.51e+04 | -1.14e+04 |
| <b>waterfront</b>    | 7489.9003  | 651.556  | 11.495  | 0.000 | 6212.797  | 8767.003  |
| <b>view</b>          | 2.385e+04  | 672.182  | 35.482  | 0.000 | 2.25e+04  | 2.52e+04  |
| <b>condition</b>     | 1.617e+04  | 690.800  | 23.406  | 0.000 | 1.48e+04  | 1.75e+04  |
| <b>grade</b>         | 5.007e+04  | 1062.780 | 47.117  | 0.000 | 4.8e+04   | 5.22e+04  |
| <b>sqft_above</b>    | 9.563e+04  | 1314.489 | 72.748  | 0.000 | 9.31e+04  | 9.82e+04  |
| <b>yr_built</b>      | -1.667e+04 | 1051.449 | -15.855 | 0.000 | -1.87e+04 | -1.46e+04 |
| <b>lat</b>           | 2.181e+04  | 4950.717 | 4.406   | 0.000 | 1.21e+04  | 3.15e+04  |
| <b>long</b>          | -9339.4918 | 3566.940 | -2.618  | 0.009 | -1.63e+04 | -2348.002 |
| <b>renovated</b>     | 5791.9948  | 641.256  | 9.032   | 0.000 | 4535.081  | 7048.909  |
| <b>has_basement</b>  | 4.717e+04  | 1636.087 | 28.833  | 0.000 | 4.4e+04   | 5.04e+04  |
| <b>zipcode_98002</b> | 6079.6824  | 7819.678 | 0.777   | 0.437 | -9247.516 | 2.14e+04  |
| <b>zipcode_98003</b> | -1.128e+04 | 6990.820 | -1.613  | 0.107 | -2.5e+04  | 2426.563  |
| <b>zipcode_98004</b> | 4.749e+05  | 1.37e+04 | 34.573  | 0.000 | 4.48e+05  | 5.02e+05  |
| <b>zipcode_98005</b> | 2.939e+05  | 1.39e+04 | 21.180  | 0.000 | 2.67e+05  | 3.21e+05  |
| <b>zipcode_98006</b> | 2.484e+05  | 1.14e+04 | 21.768  | 0.000 | 2.26e+05  | 2.71e+05  |
| <b>zipcode_98007</b> | 2.176e+05  | 1.42e+04 | 15.283  | 0.000 | 1.9e+05   | 2.46e+05  |
| <b>zipcode_98008</b> | 2.005e+05  | 1.36e+04 | 14.751  | 0.000 | 1.74e+05  | 2.27e+05  |
| <b>zipcode_98010</b> | 9.968e+04  | 1.2e+04  | 8.313   | 0.000 | 7.62e+04  | 1.23e+05  |
| <b>zipcode_98011</b> | 8.508e+04  | 1.76e+04 | 4.824   | 0.000 | 5.05e+04  | 1.2e+05   |
| <b>zipcode_98014</b> | 8.342e+04  | 1.94e+04 | 4.299   | 0.000 | 4.54e+04  | 1.21e+05  |
| <b>zipcode_98019</b> | 5.968e+04  | 1.91e+04 | 3.126   | 0.002 | 2.23e+04  | 9.71e+04  |
| <b>zipcode_98022</b> | 2.404e+04  | 1.05e+04 | 2.294   | 0.022 | 3502.443  | 4.46e+04  |
| <b>zipcode_98023</b> | -2.91e+04  | 6450.225 | -4.512  | 0.000 | -4.17e+04 | -1.65e+04 |
| <b>zipcode_98024</b> | 1.335e+05  | 1.72e+04 | 7.771   | 0.000 | 9.99e+04  | 1.67e+05  |
| <b>zipcode_98027</b> | 1.777e+05  | 1.16e+04 | 15.366  | 0.000 | 1.55e+05  | 2e+05     |
| <b>zipcode_98028</b> | 7.239e+04  | 1.72e+04 | 4.219   | 0.000 | 3.88e+04  | 1.06e+05  |
| <b>zipcode_98029</b> | 2.089e+05  | 1.32e+04 | 15.786  | 0.000 | 1.83e+05  | 2.35e+05  |
| <b>zipcode_98030</b> | 1409.0238  | 7704.114 | 0.183   | 0.855 | -1.37e+04 | 1.65e+04  |
| <b>zipcode_98031</b> | 2269.2724  | 8043.564 | 0.282   | 0.778 | -1.35e+04 | 1.8e+04   |
| <b>zipcode_98032</b> | -1.804e+04 | 9301.251 | -1.939  | 0.053 | -3.63e+04 | 195.624   |
| <b>zipcode_98033</b> | 2.701e+05  | 1.48e+04 | 18.203  | 0.000 | 2.41e+05  | 2.99e+05  |
| <b>zipcode_98034</b> | 1.25e+05   | 1.58e+04 | 7.904   | 0.000 | 9.4e+04   | 1.56e+05  |

|                      |            |          |        |       |           |          |
|----------------------|------------|----------|--------|-------|-----------|----------|
| <b>zipcode_98038</b> | 5.001e+04  | 8718.135 | 5.736  | 0.000 | 3.29e+04  | 6.71e+04 |
| <b>zipcode_98039</b> | 6.029e+05  | 3.75e+04 | 16.071 | 0.000 | 5.29e+05  | 6.76e+05 |
| <b>zipcode_98040</b> | 3.904e+05  | 1.22e+04 | 32.021 | 0.000 | 3.66e+05  | 4.14e+05 |
| <b>zipcode_98042</b> | 1.043e+04  | 7403.795 | 1.408  | 0.159 | -4085.059 | 2.49e+04 |
| <b>zipcode_98045</b> | 1.111e+05  | 1.62e+04 | 6.865  | 0.000 | 7.94e+04  | 1.43e+05 |
| <b>zipcode_98052</b> | 2.092e+05  | 1.5e+04  | 13.927 | 0.000 | 1.8e+05   | 2.39e+05 |
| <b>zipcode_98053</b> | 2.012e+05  | 1.61e+04 | 12.472 | 0.000 | 1.7e+05   | 2.33e+05 |
| <b>zipcode_98055</b> | 2.192e+04  | 8989.182 | 2.438  | 0.015 | 4297.280  | 3.95e+04 |
| <b>zipcode_98056</b> | 7.812e+04  | 9814.117 | 7.960  | 0.000 | 5.89e+04  | 9.74e+04 |
| <b>zipcode_98058</b> | 2.411e+04  | 8505.148 | 2.834  | 0.005 | 7436.354  | 4.08e+04 |
| <b>zipcode_98059</b> | 8.62e+04   | 9614.416 | 8.965  | 0.000 | 6.74e+04  | 1.05e+05 |
| <b>zipcode_98065</b> | 1.264e+05  | 1.49e+04 | 8.474  | 0.000 | 9.72e+04  | 1.56e+05 |
| <b>zipcode_98070</b> | 7.211e+04  | 1.14e+04 | 6.348  | 0.000 | 4.98e+04  | 9.44e+04 |
| <b>zipcode_98072</b> | 1.227e+05  | 1.76e+04 | 6.978  | 0.000 | 8.82e+04  | 1.57e+05 |
| <b>zipcode_98074</b> | 1.88e+05   | 1.42e+04 | 13.230 | 0.000 | 1.6e+05   | 2.16e+05 |
| <b>zipcode_98075</b> | 2.145e+05  | 1.37e+04 | 15.671 | 0.000 | 1.88e+05  | 2.41e+05 |
| <b>zipcode_98077</b> | 1.269e+05  | 1.83e+04 | 6.929  | 0.000 | 9.1e+04   | 1.63e+05 |
| <b>zipcode_98092</b> | -7660.1869 | 7003.428 | -1.094 | 0.274 | -2.14e+04 | 6067.095 |
| <b>zipcode_98102</b> | 3.42e+05   | 1.55e+04 | 22.018 | 0.000 | 3.12e+05  | 3.72e+05 |
| <b>zipcode_98103</b> | 2.507e+05  | 1.43e+04 | 17.593 | 0.000 | 2.23e+05  | 2.79e+05 |
| <b>zipcode_98105</b> | 3.089e+05  | 1.49e+04 | 20.694 | 0.000 | 2.8e+05   | 3.38e+05 |
| <b>zipcode_98106</b> | 7.308e+04  | 1.05e+04 | 6.970  | 0.000 | 5.25e+04  | 9.36e+04 |
| <b>zipcode_98107</b> | 2.482e+05  | 1.47e+04 | 16.914 | 0.000 | 2.19e+05  | 2.77e+05 |
| <b>zipcode_98108</b> | 7.488e+04  | 1.16e+04 | 6.478  | 0.000 | 5.22e+04  | 9.75e+04 |
| <b>zipcode_98109</b> | 3.433e+05  | 1.56e+04 | 21.986 | 0.000 | 3.13e+05  | 3.74e+05 |
| <b>zipcode_98112</b> | 3.688e+05  | 1.4e+04  | 26.420 | 0.000 | 3.41e+05  | 3.96e+05 |
| <b>zipcode_98115</b> | 2.535e+05  | 1.45e+04 | 17.490 | 0.000 | 2.25e+05  | 2.82e+05 |
| <b>zipcode_98116</b> | 2.366e+05  | 1.18e+04 | 20.103 | 0.000 | 2.14e+05  | 2.6e+05  |
| <b>zipcode_98117</b> | 2.387e+05  | 1.47e+04 | 16.282 | 0.000 | 2.1e+05   | 2.67e+05 |
| <b>zipcode_98118</b> | 1.236e+05  | 1.02e+04 | 12.067 | 0.000 | 1.04e+05  | 1.44e+05 |
| <b>zipcode_98119</b> | 3.321e+05  | 1.45e+04 | 22.841 | 0.000 | 3.04e+05  | 3.61e+05 |
| <b>zipcode_98122</b> | 2.459e+05  | 1.28e+04 | 19.268 | 0.000 | 2.21e+05  | 2.71e+05 |
| <b>zipcode_98125</b> | 1.291e+05  | 1.56e+04 | 8.260  | 0.000 | 9.85e+04  | 1.6e+05  |
| <b>zipcode_98126</b> | 1.386e+05  | 1.08e+04 | 12.884 | 0.000 | 1.17e+05  | 1.6e+05  |
| <b>zipcode_98133</b> | 7.658e+04  | 1.61e+04 | 4.743  | 0.000 | 4.49e+04  | 1.08e+05 |

|                                       |            |          |         |       |           |           |
|---------------------------------------|------------|----------|---------|-------|-----------|-----------|
| <b>zipcode_98136</b>                  | 2.047e+05  | 1.1e+04  | 18.531  | 0.000 | 1.83e+05  | 2.26e+05  |
| <b>zipcode_98144</b>                  | 1.918e+05  | 1.19e+04 | 16.160  | 0.000 | 1.69e+05  | 2.15e+05  |
| <b>zipcode_98146</b>                  | 7.139e+04  | 9838.519 | 7.257   | 0.000 | 5.21e+04  | 9.07e+04  |
| <b>zipcode_98148</b>                  | 3.009e+04  | 1.32e+04 | 2.275   | 0.023 | 4163.796  | 5.6e+04   |
| <b>zipcode_98155</b>                  | 6.579e+04  | 1.68e+04 | 3.917   | 0.000 | 3.29e+04  | 9.87e+04  |
| <b>zipcode_98166</b>                  | 7.008e+04  | 9010.377 | 7.778   | 0.000 | 5.24e+04  | 8.77e+04  |
| <b>zipcode_98168</b>                  | 1.88e+04   | 9482.544 | 1.983   | 0.047 | 216.730   | 3.74e+04  |
| <b>zipcode_98177</b>                  | 1.401e+05  | 1.69e+04 | 8.267   | 0.000 | 1.07e+05  | 1.73e+05  |
| <b>zipcode_98178</b>                  | 2.549e+04  | 9804.095 | 2.600   | 0.009 | 6270.930  | 4.47e+04  |
| <b>zipcode_98188</b>                  | 1.035e+04  | 9993.181 | 1.035   | 0.301 | -9241.852 | 2.99e+04  |
| <b>zipcode_98198</b>                  | 5236.4108  | 7581.752 | 0.691   | 0.490 | -9624.433 | 2.01e+04  |
| <b>zipcode_98199</b>                  | 2.926e+05  | 1.4e+04  | 20.850  | 0.000 | 2.65e+05  | 3.2e+05   |
| <b>has_larger_sqft_than_neighbors</b> | -1.719e+04 | 1718.780 | -10.003 | 0.000 | -2.06e+04 | -1.38e+04 |

**Omnibus:** 1998.403    **Durbin-Watson:** 1.987

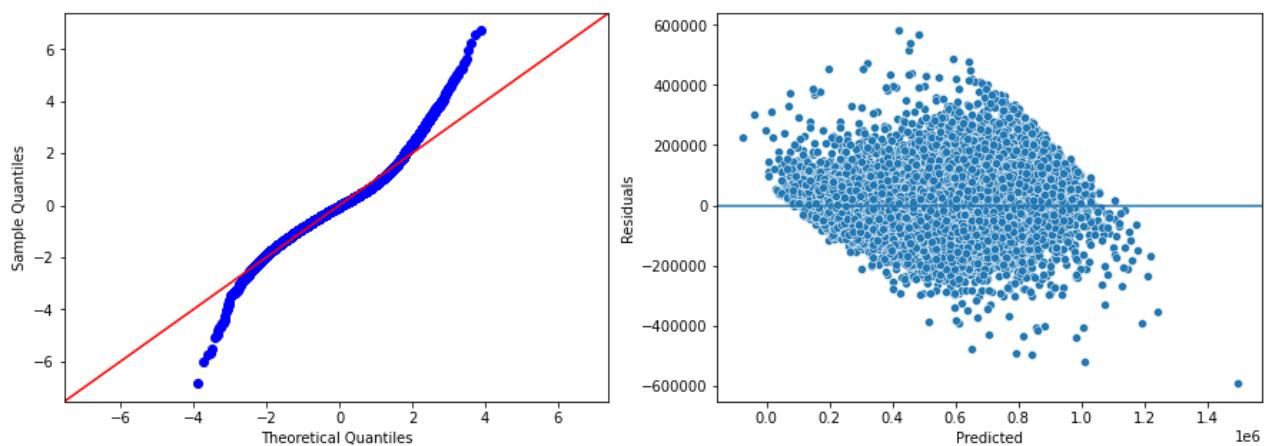
**Prob(Omnibus):** 0.000    **Jarque-Bera (JB):** 7377.212

**Skew:** 0.458    **Prob(JB):** 0.00

**Kurtosis:** 5.797    **Cond. No.** 284.

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.



There are no changes to our R-squared or p-values since we are merely scaling model 2.2 to be able to interpret effect sizes of the different parameters. The residuals also did not change since we did not change any of our model.

## iNTERPRET

## Final Scaled Coefficients

```
In [72]: coeffs = model.params.sort_values().to_frame('coeffs')
coeffs['abs'] = coeffs['coeffs'].abs()
coeffs.sort_values('abs', ascending=False, inplace=True)
```

### Non zipcode Parameters

```
In [73]: coeffs[~coeffs.index.str.startswith('zipcode')]
```

|                                       | coeffs        | abs           |
|---------------------------------------|---------------|---------------|
| <b>Intercept</b>                      | 330273.764168 | 330273.764168 |
| <b>sqft_above</b>                     | 95626.887282  | 95626.887282  |
| <b>grade</b>                          | 50074.835277  | 50074.835277  |
| <b>has_basement</b>                   | 47173.818402  | 47173.818402  |
| <b>view</b>                           | 23850.371935  | 23850.371935  |
| <b>lat</b>                            | 21810.707948  | 21810.707948  |
| <b>has_larger_sqft_than_neighbors</b> | -17192.309109 | 17192.309109  |
| <b>yr_built</b>                       | -16670.283873 | 16670.283873  |
| <b>condition</b>                      | 16169.003737  | 16169.003737  |
| <b>bathrooms</b>                      | 13774.717849  | 13774.717849  |
| <b>floors</b>                         | -13290.789002 | 13290.789002  |
| <b>sqft_lot</b>                       | 12021.914526  | 12021.914526  |
| <b>long</b>                           | -9339.491807  | 9339.491807   |
| <b>waterfront</b>                     | 7489.900314   | 7489.900314   |
| <b>renovated</b>                      | 5791.994762   | 5791.994762   |

When compared to each other, we are seeing that the top 3 parameters that affect the sales price of a home that can be changed through a renovation are the total sqft above ground, the grade of construction/finishes as well as whether the house had a finished basement or not.

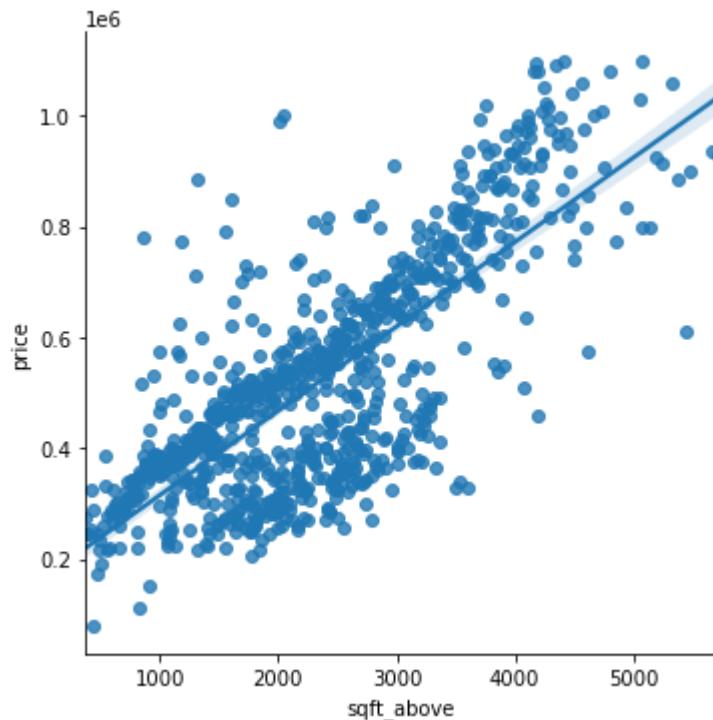
```
In [74]: df_sqft_above = df_IQR_price.loc[:, ['price', 'sqft_above']]
```

```
In [75]: df_sqft_above = df_sqft_above.groupby('sqft_above').mean().reset_index()
```

## Data Visualizations

```
In [76]: sns.lmplot(x='sqft_above', y='price', data=df_sqft_above)
```

```
Out[76]: <seaborn.axisgrid.FacetGrid at 0x24d03e1efd0>
```



```
In [77]: print(df_sqft_above['sqft_above'].min())
print(df_sqft_above['sqft_above'].max())
```

370  
5710

```
In [78]: def categorize(x):
    if (x<1000) & (x > 0):
        val = 'Up to 1000 SF'
    elif (x>=1000) & (x<2000):
        val = '1000-2000 SF'
    elif (x>=2000) & (x<3000):
        val = '2000-3000 SF'
    elif (x>=3000) & (x<4000):
        val = '3000 - 4000 SF'
    elif (x>=4000) & (x<5000):
        val = '4000 - 5000 SF'
    else:
        val = '5000+ SF'
    return val
```

```
In [79]: df_sqft_above['Category'] = df_sqft_above['sqft_above'].map(lambda x: categorize(x))
```

```
In [80]: df_sqft_above['Category'].value_counts()
```

```
Out[80]: 2000-3000 SF      272
1000-2000 SF      263
3000 - 4000 SF     151
Up to 1000 SF       78
4000 - 5000 SF      60
5000+ SF            12
Name: Category, dtype: int64
```

```
In [81]: mean_price_per_cat = df_sqft_above.groupby('Category')['price'].mean().reset_index()
```

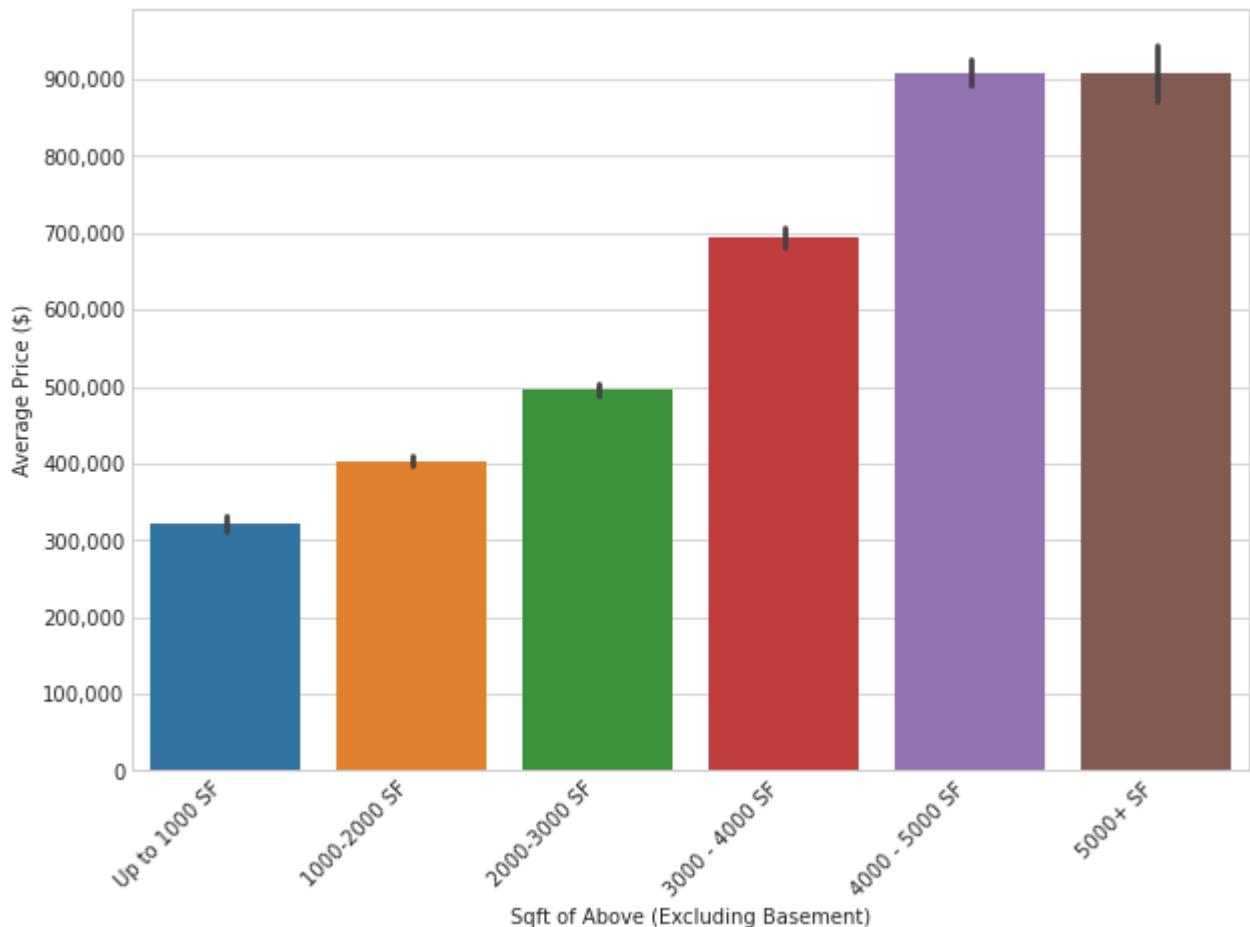
```
In [82]: mean_price_per_cat
```

Out[82]:

|   | Category       | price         |
|---|----------------|---------------|
| 0 | 1000-2000 SF   | 403132.561189 |
| 1 | 2000-3000 SF   | 496939.230583 |
| 2 | 3000 - 4000 SF | 694294.614732 |
| 3 | 4000 - 5000 SF | 908895.025397 |
| 4 | 5000+ SF       | 907491.666667 |
| 5 | Up to 1000 SF  | 321905.992281 |

In [83]:

```
from matplotlib.ticker import FuncFormatter
order = ['Up to 1000 SF', '1000-2000 SF', '2000-3000 SF', '3000 - 4000 SF', '4000 - 5000 SF', '5000+ SF']
with plt.style.context('seaborn-whitegrid'):
    fig, ax = plt.subplots(figsize=(10,7))
    sns.barplot(data = df_sqft_above, x='Category', y= 'price', order=order, ci=68)
    ax.set_xlabel('Sqft of Above (Excluding Basement)')
    ax.set_ylabel('Average Price ($)')
    ax.yaxis.set_major_formatter(FuncFormatter(lambda x, p: format(int(x), ',')))
    ax.set_xticklabels(ax.get_xticklabels(), rotation=45, ha='right')
    ax.set_yticks(range(0,1000000,100000))
```



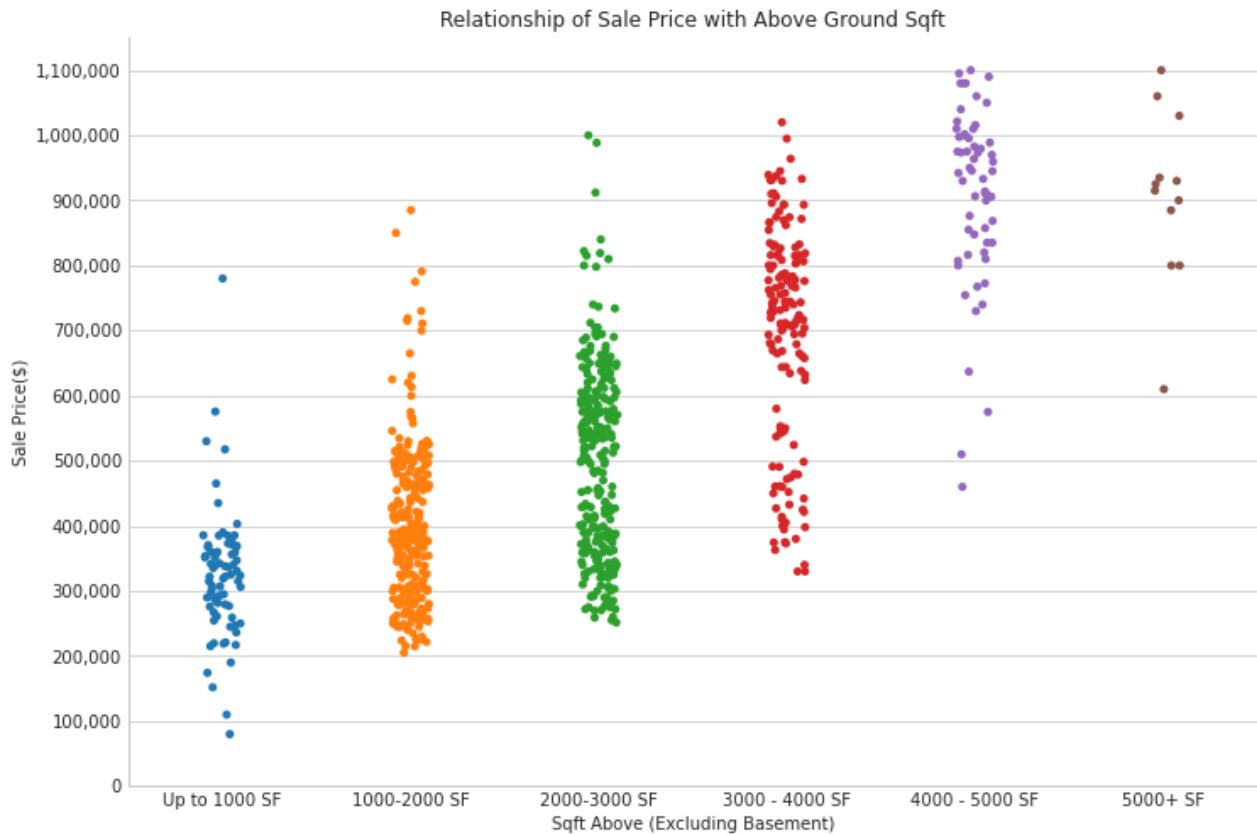
In [84]:

```
with plt.style.context('seaborn-whitegrid'):
    sns.catplot(data = df_sqft_above, x='Category', y='price', aspect=1.5, order=order)
    ax = plt.gca()
    fig = plt.gcf()
    ax.set_yticks(range(0,1200000,100000))
```

```

ax.yaxis.set_major_formatter(FuncFormatter(lambda x, p: format(int(x), ',')))
ax.set_ylabel('Sale Price($)')
fig.set_size_inches(10, 7)
ax.set_xlabel('Sqft Above (Excluding Basement)')
ax.set_title('Relationship of Sale Price with Above Ground Sqft')
#     sns.pointplot(x= 'Category', y='price', data=mean_price_per_cat, order=order)

```



As can be seen above, as the square footage of the house increased sale price of the home also tended to increase with it. Even though there is a spread of price at each category of square footage and therefore some overlaps between them, there is a clear positive trend between sale price and square footage above ground. Additionally, it should be noted that as the square footage increases the minimum sale price of the category is consistently increasing as well indicating that even if the house may not be optimal per other parameters, having a certain amount of square footage tends to make it more valuable.

```
In [85]: df_IQR_price['has_basement'] = df_IQR_price['has_basement'].astype(bool)
```

```
<ipython-input-85-8ecde1b6317e>:1: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead
```

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

```
df_IQR_price['has_basement'] = df_IQR_price['has_basement'].astype(bool)
```

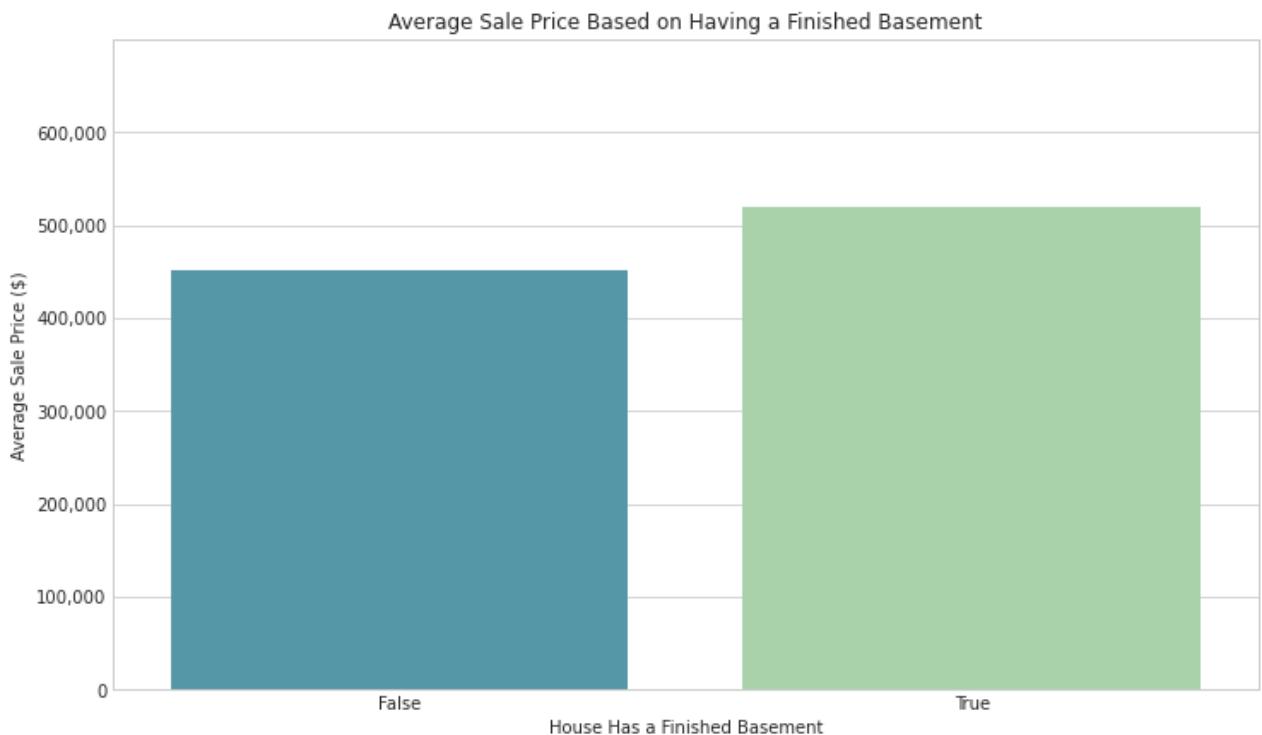
```
In [86]: df_basement = df_IQR_price.loc[:,['price','has_basement']]
df_basement_price = df_basement.groupby('has_basement')['price'].mean()
```

```
In [87]: df_basement_price = df_basement_price.to_frame().reset_index()
```

```
In [88]: df_basement_price
```

| Out[88]: | has_basement | price         |
|----------|--------------|---------------|
| 0        | False        | 451003.365626 |
| 1        | True         | 520135.677810 |

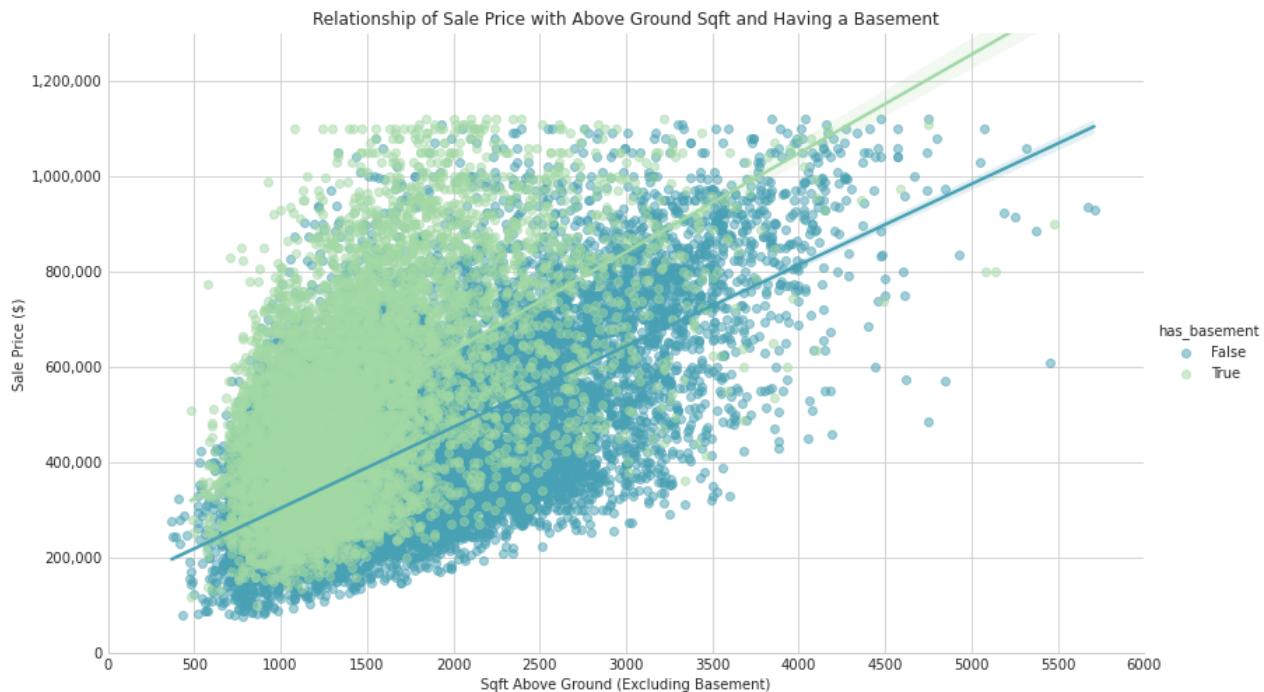
```
In [89]: with plt.style.context('seaborn-whitegrid'):
    fig, ax = plt.subplots(figsize=(12,7))
    sns.barplot(data = df_basement_price, x='has_basement', y= 'price', palette=sns.color_palette("Spectral"))
    ax.set_xlabel('House Has a Finished Basement')
    ax.set_ylabel('Average Sale Price ($)')
    ax.set_title('Average Sale Price Based on Having a Finished Basement')
    ax.set_yticks(range(0,700000,100000))
    ax.yaxis.set_major_formatter(FuncFormatter(lambda x, p: format(int(x), ',')))
    ax.set_xticks([0,700000,1000000]);
```



Exploring the relationship between having a finished basement and the sale price, we see that homes that had finished basements had a higher average sale price of approximately 520,000 dollars compared to homes that did not at approximately 451,000 dollars.

```
In [90]: from matplotlib.ticker import FuncFormatter
with plt.style.context('seaborn-whitegrid'):

    sns.lmplot(x='sqft_above', y='price', data=df_IQR_price, hue='has_basement',
                aspect=2, scatter_kws=dict(alpha=0.5), palette=sns.color_palette("Spectral"))
    ax = plt.gca()
    ax.set_xlabel('Sqft Above Ground (Excluding Basement)')
    ax.set_ylabel('Sale Price ($)')
    ax.set_title("Relationship of Sale Price with Above Ground Sqft and Having a Basement")
    ax.yaxis.set_major_formatter(FuncFormatter(lambda x, p: format(int(x), ',')))
    ax.xaxis.set_ticks(range(0,6500,500))
    ax.set_yticks([0, 1300000])
    plt.gcf().set_size_inches(12,7);
```



When we look at the relationship between square footage above ground and the home's sale price again, but add in a second parameter to define whether the house had a basement or not, the houses with a basement have a slightly more positive relationship with the sale price. So if a house had a basement, it tended to have a slightly higher price than a comparable home. This is visible from the difference of slopes between the two lines shown above where the green line is diverging from the blue line in a positive way.

```
In [91]: df_grade = df_IQR_price.loc[:,['price','grade']]
df_grade
```

```
Out[91]:      price  grade
```

|       | price    | grade |
|-------|----------|-------|
| 0     | 221900.0 | 7     |
| 1     | 538000.0 | 7     |
| 2     | 180000.0 | 6     |
| 3     | 604000.0 | 7     |
| 4     | 510000.0 | 8     |
| ...   | ...      | ...   |
| 21592 | 360000.0 | 8     |
| 21593 | 400000.0 | 8     |
| 21594 | 402101.0 | 7     |
| 21595 | 400000.0 | 8     |
| 21596 | 325000.0 | 7     |

20439 rows × 2 columns

```
In [92]: def categorize_grade(x):
```

```

if (x<7):
    val = 'Below Average'
elif (x<=8) & (x>6):
    val = 'Average'
else:
    val = 'Above Average'

return val

```

In [93]: df\_grade['Category'] = df\_grade['grade'].map(lambda x: categorize\_grade(x))

In [94]: df\_grade['Category'].value\_counts()

Out[94]:

| Category      | Count |
|---------------|-------|
| Average       | 14905 |
| Above Average | 3227  |
| Below Average | 2307  |

Name: Category, dtype: int64

In [95]: df\_grade  
df\_grade\_sqft = pd.concat([df\_grade, pd.DataFrame(df\_IQR\_price['sqft\_above'])], axis=1)  
df\_grade\_sqft

Out[95]:

|       | price    | grade | Category      | sqft_above |
|-------|----------|-------|---------------|------------|
| 0     | 221900.0 | 7     | Average       | 1180       |
| 1     | 538000.0 | 7     | Average       | 2170       |
| 2     | 180000.0 | 6     | Below Average | 770        |
| 3     | 604000.0 | 7     | Average       | 1050       |
| 4     | 510000.0 | 8     | Average       | 1680       |
| ...   | ...      | ...   | ...           | ...        |
| 21592 | 360000.0 | 8     | Average       | 1530       |
| 21593 | 400000.0 | 8     | Average       | 2310       |
| 21594 | 402101.0 | 7     | Average       | 1020       |
| 21595 | 400000.0 | 8     | Average       | 1600       |
| 21596 | 325000.0 | 7     | Average       | 1020       |

20439 rows × 4 columns

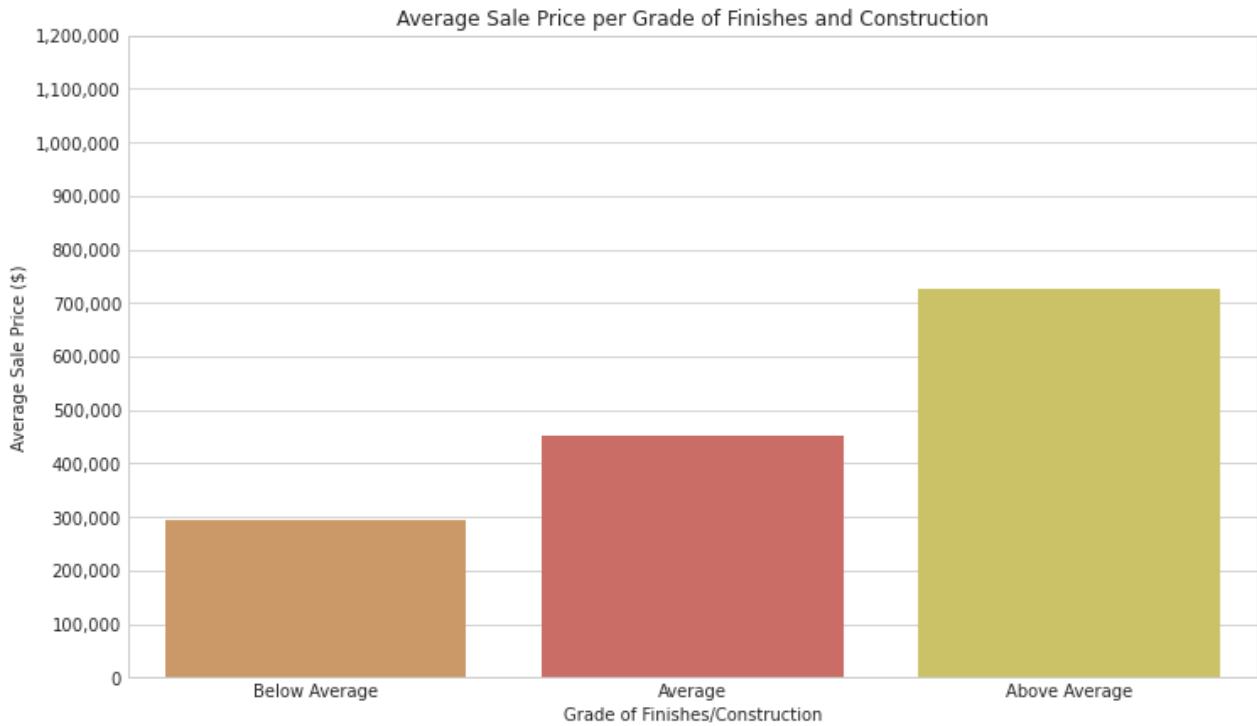
In [96]: mean\_price\_per\_grade = df\_grade.groupby('Category')['price'].mean().reset\_index()  
mean\_price\_per\_grade

Out[96]:

| Category        | price         |
|-----------------|---------------|
| 0 Above Average | 726319.069724 |
| 1 Average       | 450805.678497 |
| 2 Below Average | 294856.879931 |

In [97]: with plt.style.context('seaborn-whitegrid'):
order=['Below Average', 'Average', 'Above Average']

```
p = sns.color_palette("hls", 14)
fig, ax = plt.subplots(figsize=(12,7))
sns.barplot(data = mean_price_per_grade, x='Category', y= 'price', order=order,
            palette=[p[1], p[0], p[2]])
ax.set_xlabel('Grade of Finishes/Construction')
ax.set_ylabel('Average Sale Price ($)')
ax.set_title('Average Sale Price per Grade of Finishes and Construction')
ax.set_ylim(0,1200000)
ax.yaxis.set_major_formatter(FuncFormatter(lambda x, p: format(int(x), ',')))
ax.set_yticks(range(0,1300000,100000))
```



Another interesting, but also expected relationship is between the grade of finishes and the sale price. Houses that had higher grades of finishes and a better construction quality sold for higher prices.

```
In [98]: from matplotlib.ticker import FuncFormatter
with plt.style.context('seaborn-whitegrid'):

    sns.lmplot(x='sqft_above', y='price', data=df_grade_sqft, hue='Category',
                aspect=2, scatter_kws=dict(alpha=0.7), palette=sns.color_palette("hls",
                ax = plt.gca()
                ax.set_xlabel('Sqft Above Ground (Excluding Basement)')
                ax.set_ylabel('Sale Price ($)')
                ax.set_title("Relationship of Sale Price with Above Ground Sqft and Grade")
                ax.yaxis.set_major_formatter(FuncFormatter(lambda x, p: format(int(x), ',')))
                ax.xaxis.set_ticks(range(0,6500,500))
                ax.set_ylim(0, 1300000)
                plt.gcf().set_size_inches(12,7);
```



When grade is plotted with square footage above ground, the relationship shown in the previous visual becomes even more apparent. As can be seen from the different colored data points, above average homes tended to have a higher sale price. The regression lines for average and above average have a higher positive slope compared to the below average homes, which means that as the square footage of a home increases, the higher graded homes will tend to have higher prices.

## CONCLUSIONS & RECOMMENDATIONS

Even though renovations are usually a lot of effort and stressful to lots of homeowners, they may help increase the property's value. To sum up, our analysis for King County, Washington showed the following:

- Increasing the square footage above ground tends to increase the house's value.
- Focusing on the grade of finishes and the quality of construction as a whole tends to pay dividends when it comes to selling the house.
- Having a basement is the third most effective parameter in increasing a home's sale price.

Our first recommendation based on our findings above would be for the homeowner to renovate their house and add livable above ground square footage to the property. This could be in the shape of an added extension to the home or a simpler approach of finishing the attic space.

Secondly, we would advise the homeowner to also focus their renovation efforts on the finishes and quality of the materials that they would be choosing to use. This could include components such as kitchen countertop material, cabinetry design, bathroom fixtures, general lighting design and fixture choices inside and outside the house, flooring.

Lastly, we would propose that the homeowner add a basement to their home. It should be noted that since we do not know what the construction costs of this undertaking would be like versus the gain in the home's sale price, it is difficult to predict the outcome (refer to the Limitations & Next

Steps section for more information on this). However, based on our model, we found that having a basement for homes was the third most effective parameter. If the home already has a basement, but the basement is unfinished, our recommendation would be for the homeowner to make the basement an occupiable space by finishing it.

## Limitations & Next Steps

Given more time and information about what the homeowner's renovation budget would be, we would have wanted to analyze whether these top 3 parameters would truly be the most effective in bringing a net value increase since a renovation such as adding a basement to a home would be very costly and may not end up returning a net value increase. Additionally, the construction costs in the state of Washington may be higher than other states due to factors such as permitting, material costs, logistical challenges etc. which may effect the net value increase as well. Furthermore, having information about whether the homeowner is thinking about living in the renovated house or renting it out would allow us to fine tune our analysis and bring more valuable insight.

## Appendix A: Effect of Zipcodes on Home Sale Price

```
In [99]: coeffs[coeffs.index.str.startswith('zipcode')]
```

|                      | coeffs        | abs           |
|----------------------|---------------|---------------|
| <b>zipcode_98039</b> | 602853.692039 | 602853.692039 |
| <b>zipcode_98004</b> | 474874.562216 | 474874.562216 |
| <b>zipcode_98040</b> | 390372.947380 | 390372.947380 |
| <b>zipcode_98112</b> | 368770.378411 | 368770.378411 |
| <b>zipcode_98109</b> | 343272.529402 | 343272.529402 |
| ...                  | ...           | ...           |
| <b>zipcode_98092</b> | -7660.186920  | 7660.186920   |
| <b>zipcode_98002</b> | 6079.682440   | 6079.682440   |
| <b>zipcode_98198</b> | 5236.410829   | 5236.410829   |
| <b>zipcode_98031</b> | 2269.272420   | 2269.272420   |
| <b>zipcode_98030</b> | 1409.023787   | 1409.023787   |

69 rows × 2 columns

```
In [100...]: c = coeffs.reset_index()
c.head()
```

|          | index         | coeffs        | abs           |
|----------|---------------|---------------|---------------|
| <b>0</b> | zipcode_98039 | 602853.692039 | 602853.692039 |

|   | index         | coeffs        | abs           |
|---|---------------|---------------|---------------|
| 1 | zipcode_98004 | 474874.562216 | 474874.562216 |
| 2 | zipcode_98040 | 390372.947380 | 390372.947380 |
| 3 | zipcode_98112 | 368770.378411 | 368770.378411 |
| 4 | zipcode_98109 | 343272.529402 | 343272.529402 |

```
In [101... #database from https://www.unitedstateszipcodes.org/zip-code-database/
df_zipcode = pd.read_csv('reference/zip_code_database.csv')
```

```
In [102... df_zipcode[df_zipcode['zip']==98004]
```

|       | zip   | type     | decommissioned | primary_city | acceptable_cities                                 | unacceptable_cities | state              | c  |
|-------|-------|----------|----------------|--------------|---|---------------------|--------------------|----|
| 41628 | 98004 | STANDARD | 0              | Bellevue     | Beaux Arts, Clyde Hill, Hunts Point, Yarrow Point |                     | Beaux Arts Village | WA |

```
In [103... zipcodes = {}
for val in c['index']:
    zipcodes[val[8:13]] = None
```

```
In [104... for value in list(zipcodes.keys()):
    if value.startswith('9'):
        continue
    else:
        del zipcodes[value]
```

```
In [105... df_zipcode['zip'] = df_zipcode['zip'].map(lambda x: str(x))
```

```
In [106... df_zipcode[df_zipcode['zip'] == '98004']
```

|       | zip   | type     | decommissioned | primary_city | acceptable_cities                                 | unacceptable_cities | state              | c  |
|-------|-------|----------|----------------|--------------|---|---------------------|--------------------|----|
| 41628 | 98004 | STANDARD | 0              | Bellevue     | Beaux Arts, Clyde Hill, Hunts Point, Yarrow Point |                     | Beaux Arts Village | WA |

```
In [107... idx = []
for key in zipcodes.keys():
    for idx, cell in enumerate(df_zipcode['zip']):
        if cell == key:
            zipcodes[key] = df_zipcode.loc[idx, 'primary_city']
```

```
In [108... #top 5 cities
print("Top 5 cities:", list(zipcodes.values())[0:5])

#top 5 zipcodes
print("Top 5 zipcodes:", list(zipcodes.keys())[0:5])
```

Top 5 cities: ['Medina', 'Bellevue', 'Mercer Island', 'Seattle', 'Seattle']

Top 5 zipcodes: ['98039', '98004', '98040', '98112', '98109']

These are the top 5 cities that had the highest effect on a home's sales price. Since the homeowner can not change their zipcodes with a renovation, we are only including this information as reference.